

United States Naval Postgraduate School



THESIS

A METHOD OF CLUSTER ANALYSIS

by

William Michael Cima

June 1970

Thesis
C4789

This document has been approved for public
release and sale; its distribution is unlimited.

United States Naval Postgraduate School



THESIS

A METHOD OF CLUSTER ANALYSIS

by

William Michael Cima

June 1970

This document has been approved for public
release and sale; its distribution is unlimited.

A Method of Cluster Analysis

by

William Michael Cima
Lieutenant (junior grade), United States Navy
B.S., United States Naval Academy, 1969

Submitted in partial fulfillment of the
requirements for the degree of

MASTER OF SCIENCE IN OPERATIONS RESEARCH

from the

NAVAL POSTGRADUATE SCHOOL
June 1970

ABSTRACT

A method of cluster analysis is presented in which points in n -dimensional space are analyzed through a subdivisive procedure. The points are orthogonally projected onto that line which maximizes their variance and the resulting point distribution is then analyzed with the use of a histogram. Wherever possible, divisions between conglomerates of points are made and each separate clump is subsequently analyzed. Ultimately adjacent groups are combined and analyzed through an analogous technique in an effort to re-unite any points which may have inadvertently deviated from the group with which they truly associate. The method is later refined to allow the detection of groups in several point dispersions which would have appeared as a single conglomeration under the original method. An example is given to illustrate the applicability of the procedure.

TABLE OF CONTENTS

I.	INTRODUCTION-----	5
	A. EXPLANATION OF THE PROBLEM -----	5
	B. GENERAL -----	6
	C. HISTORY -----	7
II.	PROCEDURE -----	11
	A. BACKGROUND -----	11
	B. THE APPROPRIATE LINE -----	13
	C. UNDERLYING THEORY -----	18
	D. THE ANALYSIS OF PROJECTED POINTS -----	21
	E. CORRECTING THE ERROR -----	26
	F. REFINEMENT OF THE PROCEDURE -----	31
	G. INTERPRETATIONS -----	37
	H. CAUTION -----	42
III.	EXAMPLE -----	44
	A. EXPLANATION OF THE DATA -----	44
	B. PROCEDURE -----	45
	C. CONCLUSION -----	49
	APPENDIX A: HISTOGRAM ANALYSIS -----	50
	COMPUTER PROGRAM -----	61
	LIST OF REFERENCES -----	72
	INITIAL DISTRIBUTION LIST -----	74
	FORM DD 1473 -----	75

ACKNOWLEDGEMENTS

The author would like to express sincere thanks to Professor G. T. Howard for making available his time and talents which attributed much to the more technical aspects of this work. Special and greatly deserved thanks are extended to Professor Gary A. Tuck who did much to illuminate the subject matter. Through his authentic interest, patience and knowledge, he successfully led and directed the author toward the culmination of this work.

I. INTRODUCTION

A. EXPLANATION OF THE PROBLEM

It is frequently necessary to classify or distinguish between events or objects. The course of one's everyday life involves numerous automatic and casual classifications. For example, without an acute conscious awareness it is possible to discriminate between man and woman, bees and birds, or various colors. In effect, an entity is characterized through the outcome of the observations of several classification characteristics or variates. In cases where the classification is less obvious, each of these n classification variates may be treated as a dimension in n -dimensional space. The individual specimens on which the observations are being made may then be conceived of as points or n -dimensional vectors plotted on a multi-dimensional coordinate system whose axes are the scales of the specimen characteristics. Thus, every subject can be represented by one point in space.

Cluster analysis is a method of estimating the true groupings of these points. If one were observing characteristics of a sample from k different populations, $p_1, p_2, \dots, p_{i-1}, p_i, p_{i+1}, \dots, p_k$, he would then want to divide n -dimensional space into k mutually exclusive and collectively exhaustive regions, $r_1, r_2, \dots, r_{i-1}, r_i, r_{i+1}, \dots, r_k$, with the rule or procedure of assigning an individual to p_i if he belongs to r_i . In other words, cluster analysis

circumscribes the k areas of high point density in hyperspace with multi-dimensional boundaries which define divisions between the k populations.

A procedure which could accomplish this has wide applicability. For example, the dispersion of human chromosomes as they appear on a photograph of mitosis involves a two-dimensional array of points; the dispersion of stars involves a three-dimensional array: bacterial classification could involve an n -dimensional array as would any taxonomical problem in which n characters of m species are treated as variates and represented by m points in n -dimensional space. Information on the true grouping of the species could be obtained from the point grouping as seen through cluster analysis. Thus, most investigations into the classifications of objects or dispersion of things require the elucidation of spatial groupings.

B. GENERAL

There are two distinct methods of attacking the cluster analysis problem. The first of these is the agglomerative method which starts with a single point as the "group" and adds to it others which satisfy certain criteria. A group is formed when no new points can be added to the existing conglomeration; another starting point is then chosen and the process iterates to determine the second group. This technique is continued until all points are allocated to one of the k groups. Thus, groups are built up from the individual points and in this manner n -dimensional space is

divided into k mutually exclusive and collectively exhaustive regions. The other method of approach employs subdivisive techniques which arrive at an ultimate partition of n -dimensional space through an initial partition, or possibly several initial partitions, of n -dimensional space and subsequent divisions of each of those partial spaces.

C. HISTORY

The area of cluster analysis has received ever increasing attention over the last four decades by those interested in classification methods. Four of the more renowned procedures are briefly reviewed; if it is desired to learn more concerning any of these methods the list of references will be helpful in directing the reader to sources of greater detail. In 1933, Hotelling devised a principle component analysis method of grouping points which involved the successive elimination of dimensions. Points are projected orthogonally from a multi-dimensional space into a space of fewer dimensions which retains "maximum information." Said differently, those characteristics having the least amount of variability among their observations are eliminated and the space reduced. This process is continued until the points are finally in an observable space.

The weighted mean pair method developed by Sokal and Michener (1958) and altered by Rogers (1959) has been applied to entomological problems. The procedure, again updated in 1963 by Sokal and Sneath, involves operations performed with an initially constructed m by m symmetric

matrix, where m is the number of points. The resemblance between two specimens is the proportion of the number of characteristics measured which they have in common. Thus, the larger this similarity ratio, the more alike are the specimens. Then a distance, d_{ij} , between any two points, c_i and c_j , $i, j = 1, 2, \dots, m$, is defined as $d_{ij} = -\log s_{ij}$, where s_{ij} is the similarity ratio between points c_i and c_j . It is interesting to notice that as s_{ij} approaches 1, d_{ij} approaches 0. Thus, the symmetric similarity matrix (d_{ij}) is in effect a mileage chart in semi-metric space, showing distance between any two points c_i and c_j . Then an overall number, H_i , is computed for each point c_i , where $H_i = \sum_j d_{ij}$, $j = 1, 2, \dots, m$, and $i = j$. The point possessing the smallest H_i value is designated as the prime node and becomes the first member of the group. Without loss of generality it can be labeled c_1 . The point closest to c_1 , c_i , is then added to the group if its similarity coefficient satisfies certain specified criteria. With this, the first and i^{th} rows and columns of the matrix are deleted, one new row and column representing the combination of the similarity coefficients of c_1 and c_i added, and the process repeated for the $n-1$ order matrix. This method continues until no new points have similarity coefficients which satisfy the group entry criterions. Then a smallest H_i value is computed for those points which were not admitted to the group and the process iterates to form multi-groupings until all points are allocated to a group.

In 1959, Williams and Lambert concerned themselves with the case of an $n \times m$ data matrix X , where n is the number of points and m is the number of variates, consisting entirely of presence or absence data (1 or 0 respectively). Originally they applied their method to taxonomic problems in ecology, where the variates were the different plant characteristics present or absent in m individuals. The m points are divided into two subsets on the basis of the variate k which "best" separates them (in a well-defined sense): One set being those that contain k , the other being those that do not. If it is feasible to divide the groups with respect to any of the remaining variates, the appropriate dimension is chosen and the process continued.

Edwards and Cavalli-Sforza (1965) developed an accurate but highly tedious subdivisive procedure. Noting that the best division between clusters would be that which resulted in the two clusters being as dense as possible, they proceed to separate the points into two groups through every possible division of points. They then choose that division which maximizes the between group variance and minimizes the within group variance. The procedure is repeated for each group with a weighting factor, the associated between clusters sum of squares, to describe the importance of each division. Note that the drawback to this method lies in the computational labor of examining all splits. The authors admit that $(n-1)^2 2^{n-1}$ seconds are required on a computer

with a five microsecond access time; that is, twenty-one points require 100 hours and forty-one points require 54,000 years.

Presented in the following pages is a method of cluster analysis. An example is given to show how the procedure is of assistance in analyzing problems of a taxonomical nature.

II. PROCEDURE

A. BACKGROUND

The procedure presented in this paper employs subdivisive techniques. The method involves the orthogonal projection of points in n -dimensional space onto a one-dimensional line and an analysis of the resulting point distribution. In order to facilitate the group detection process, it seems desirable to determine that line which provides maximum separation between the group projections. The difficulty in doing this emanates from not knowing which points associate to form groups. A second impediment is that the meaning of maximum separation between groups is not clear. The concept of maximizing separation between groups is subject to several different interpretations, two of which follow: the line providing maximum separations between groups is obtained when (1) the sum of the distance between the closest points of all adjacent groups is a maximum or (2) any two adjacent groups have a maximum separation, that is, the separation between the projections of two adjacent groups is greater than the separation between the projections of any other adjacent groups on any line.

Figure 1 may be helpful in clarifying this distinction:

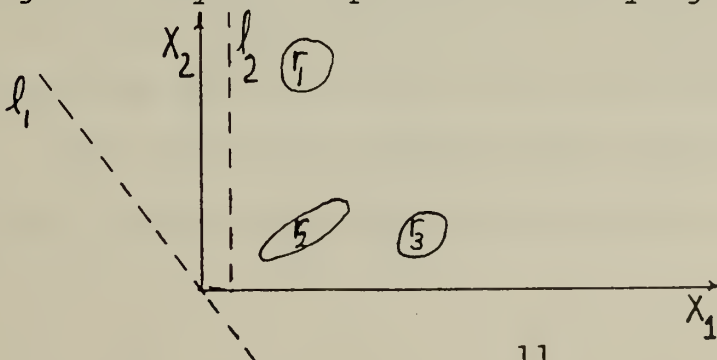


Figure 1

The line on which the sum of the distances between the closest points of all adjacent group projections is a maximum is ℓ_1 . That is, the sum of the distance between groups one and two plus that between groups two and three is greatest when projected on line ℓ_1 . It thus corresponds to the first interpretation. The maximum separation between the projections of two groups, r_1 and r_2 is provided by ℓ_2 . There is no line on which the projections of any two adjacent groups are separated by a distance greater than that between groups r_1 and r_2 on line ℓ_2 . Thus, ℓ_2 provides maximum separation with respect to the second interpretation.

Inherent to subdivisive methods is the possibility of erring by assigning a point to a group with which it does not conglomerate. It is, therefore, desirable to make divisions between groups in a manner which minimizes the probability of introducing such error. Thus, in every case it is advantageous to work with lines analogous to ℓ_2 because it provides a maximum amount of separation between two high density point areas which are to be divided into two groups. In this manner, it reduces the probability of erring by assigning points to the wrong group. However, certain point dispersions which result in well-defined groups enhance the utility of lines analogous to ℓ_1 , which has the advantage of dictating more divisions per iteration.

One line which maximizes the separation between groups under the first interpretation is that line which maximizes

the between group variance and minimizes the within group variance. However, inadequate knowledge concerning those vectors which conglomerate or which points associate to form groups makes this method non-functional computationally; in this form it involves the analysis of all possible divisions of points into groups and the subsequent projection of each combination onto the most suitable line. The computational labor of such an analysis rivals that of the method of Edwards and Cavalli-Sforza mentioned earlier.

B. THE APPROPRIATE LINE

The line which is used is not as discriminating as that which maximizes the between group variance and minimizes the within group variance, but is computationally more accessible. The basic analytical tool is chosen to be that one-dimensional ray on which the variance of the orthogonally projected points is a maximum. Several point distributions can be formed for which that line which maximizes the variance of the projected points does not coincide with or even lie near that line which maximizes the separation between groups. However, it does appear to be a reasonable line in that in most cases it provides a point distribution which is very auspicious to group detection in later analysis. The problems arising when this line is not favorable are reckoned with later. The problem of finding a line which maximizes the variance of the orthogonally projected points can be formulated as a non-linear programming problem. The objective function is to maximize the variance

of the projected points on a line given by the direction of u , a unit vector.

The problem is:

$$\text{Max}_u S^2 = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})^2$$

$$\text{Subject to: } u^t u = 1,$$

where:

$$S^2 = \text{sample variance}$$

$$m = \text{number of points}$$

$$x_i = \sum_{j=1}^n c_{ij} u_j$$

where:

$$n = \text{the number of dimensions or variables}$$

$$u^t = (u_1, u_2, \dots, u_n) = \text{unit vector}$$

$$c_{ij} = (c_{i1}, c_{i2}, \dots, c_{in}) = \text{coordinates of the } i^{\text{th}} \text{ point in } n \text{ dimensions}$$

$$\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n (c_{ij} u_j).$$

The problem now becomes:

$$\text{Max } S^2 = \frac{1}{m-1} \sum_{i=1}^m \left[\left(\sum_{j=1}^n c_{ij} u_j \right) - \left(\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n c_{ij} u_j \right) \right]^2$$

subject to $u^t u = 1$.

The Lagrangian is then:

$$F(u_1, \dots, u_n, \lambda) = S^2 - \lambda(u^t u - 1)$$

$$\text{but } u^t u = \sum_{j=1}^n u_j^2,$$

which implies

$$F(u_1, \dots, u_n, \lambda) = s^2 - \lambda \left(\sum_{j=1}^n u_j^2 - 1 \right),$$

and substituting in the expression for s^2 gives

$$F(u_1, \dots, u_n, \lambda) = \frac{1}{m-1} \sum_{i=1}^m \left[\left(\sum_{j=1}^n c_{ij} u_j \right) - \left(\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n c_{ij} u_j \right) \right]^2 - \lambda \left(\sum_{j=1}^n u_j^2 - 1 \right).$$

Setting the partial derivative of the Lagrangian with respect to each variable equal to zero gives necessary conditions from which the optimal vector, u , can be found.

For the k^{th} component of u :

$$0 = \frac{\partial F}{\partial u_k} = \frac{2}{m-1} \sum_{i=1}^m \left[\left(\sum_{j=1}^n c_{ij} u_j \right) - \left(\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n c_{ij} u_j \right) \right] \left[c_{ik} - \frac{1}{m} \sum_{i=1}^m c_{ik} \right] - 2 \lambda u_k.$$

Letting $c_{ik} - \frac{1}{m} \sum_{i=1}^m c_{ik} = B_{ik}$ and rewriting gives:

$$0 = \frac{2}{m-1} \sum_{i=1}^m \left[\left(\sum_{j=1}^n c_{ij} u_j \right) (B_{ik}) - \left(\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n c_{ij} u_j \right) (B_{ik}) \right] - 2 \lambda u_k.$$

Distributing the first summation over the terms:

$$0 = \frac{2}{m-1} \sum_{i=1}^m \left[\sum_{j=1}^n c_{ij} u_j \right] B_{ik} - \frac{2}{m(m-1)} \sum_{i=1}^m \left[\sum_{i=1}^m \sum_{j=1}^n c_{ij} u_j \right] B_{ik} - 2 \lambda u_k.$$

Rearranging within each term:

$$0 = \sum_{i=1}^m \frac{2 B_{ik}}{m-1} \sum_{j=1}^n c_{ij} u_j - \left[\sum_{i=1}^m \frac{2 B_{ik}}{m(m-1)} \right] \left[\sum_{i=1}^m \sum_{j=1}^n c_{ij} u_j \right] - 2 \lambda u_k.$$

The second term is identically zero since

$$\sum_{i=1}^m B_{ik} = \sum_{i=1}^m \left(c_{ik} - \frac{1}{m} \sum_{i=1}^m c_{ik} \right) = 0 ,$$

which shows that

$$\sum_{i=1}^m \frac{2 B_{ik}}{m(m-1)} \left[\overline{mx} \right] = 0 .$$

Therefore,

$$\frac{\partial F}{\partial u_k} = 0 = \sum_{i=1}^m \frac{2 B_{ik}}{m-1} \sum_{j=1}^n c_{ij} u_j - 2 \lambda u_k .$$

Setting $\frac{2 B_{ik}}{m-1} = D_{ik}$ yields

$$(1) \quad 0 = \sum_{i=1}^m D_{ik} \sum_{j=1}^n c_{ij} u_j - 2 \lambda u_k , \quad h = 1, \dots, n, \text{ as}$$

a necessary condition. Differentiating the Lagrangian with respect to the Lagrange multiplier, λ , and setting the derivative equal to zero yields the constraint

$$(2) \quad u_1^2 + \dots + u_n^2 = 1 .$$

Thus, u is found through the solution of a system of $n+1$ equations in $n+1$ unknowns, u_1, u_2, \dots, u_n . There are n equations of the form (1) and there is one equation of form (2). It is important to notice that setting the first partial derivatives of the Lagrangian equal to zero is a necessary but not a sufficient condition for a maximum. In fact, on rare occasions when the rank of the Jacobian of the constraint equations is less than m , the number of

constraints, these conditions are not even necessary. That is, for a set of constraints possessing a singular Jacobian, there may exist solutions which maximize or minimize the objective functions which will not be found using the first partial derivative of the Lagrangian technique. This will be ignored here, however. Since there is no guarantee that a solution to the above system of equations maximizes the objective function, it is imperative that all solutions be obtained. The unit vector which maximizes the variance of the projected points will be one of those solutions and can be identified as the maximum through its substitution into the objective function.

C. UNDERLYING THEORY

Before delving further into the details of the procedure, one assumption can be made from which two important observations follow; it is hypothesized that the n-dimensional hypervolume spanned by the points of any group is nearly a hyperellipsoid whose lengths along the axes are real numbers greater than or equal to zero. This assumption seems consistent with expectations since a hyperellipsoid is the shape of the hypervolume that the points will most likely span. Furthermore, it imposes virtually no constraint since it permits the spatial dispersion of points of any group to generate a multi-dimensional hypervolume of many forms, from a hyperball (solid hypersphere) to a straight line segment.

The assumption that the hypervolume spanned by the points is a hyperellipsoid results in the first observation --that a curve-fitted histogram plot of the distribution of points of any group projected on a line has a bell-like structure similar to that of the normal density; furthermore, the effect is especially inherent to groups of high point density. This phenomenon can be reasoned as follows:

l_1 , will represent the line onto which the points are to be orthogonally projected and the two-dimensional case will be considered. A sufficient number of equidistant parallel lines, w_1, w_2, \dots, w_j , can be constructed which completely bound and divide the entire group. The value of j would depend upon the distance between lines. In Figure 2 below, it can be observed that the expected number of points of the

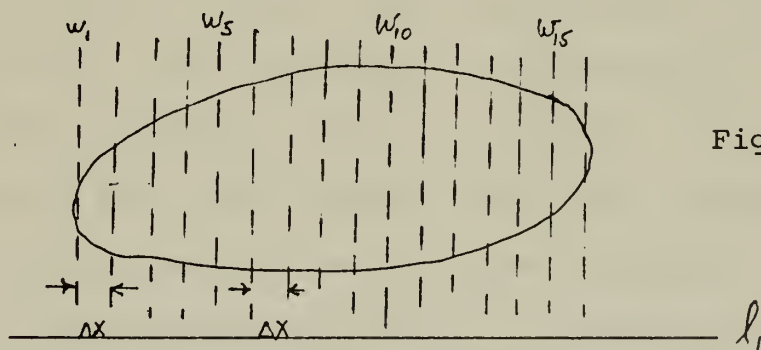


Figure 2

group contained between w_8 and w_9 is greater than that contained between w_3 and w_4 , which in turn is greater than that between w_1 and w_2 . Hence, the curve-fitted histogram can be expected to exhibit a bell-like structure. Figure 3 indicates that the orientation of the group with respect to the line is free to vary without altering this result. Now generalizing this two-dimensional analysis for the n -dimensional case, l_1 , can again represent

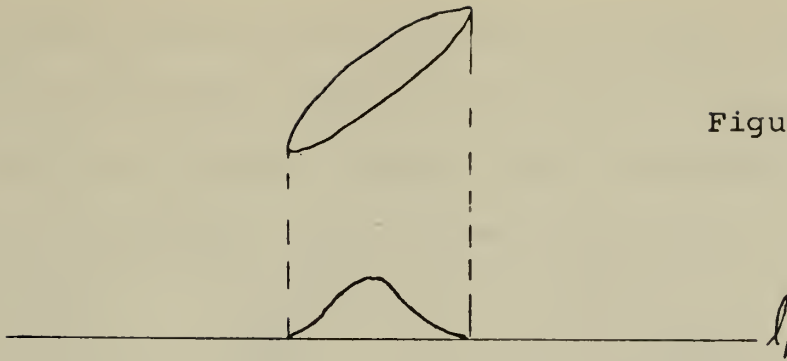


Figure 3

the line onto which the points are to be projected. Then, as before, it is possible to construct a sufficient number of equidistant parallel $(n-1)$ -dimensional hyperplanes, w_1, w_2, \dots, w_j , to completely bound and divide the hyper-ellipsoid. Then, intuitively, the expected number of points contained between two adjacent hyperplanes passing through a "thick" portion of the ellipsoid (near its geometric center) is greater than that contained between two hyperplanes in a "thinner" portion. Hence, the curve-fitted histogram would possess a bell-like structure, an effect which is more salient for increasing density of points. An interesting extension of this observation is that the normal-like shape becomes more pronounced as the hyper-ellipsoid tends to a hyperball and less pronounced as it elongates; in fact, the bell-like structure is not discernible for a hyperellipsoid in the form of a straight line. The concept of ascertaining information concerning the group shape through the analysis of point distributions on the line will later be shown to be of value.

The second observation follows from the first: a histogram of the projections of two or more groups superimposed on the line will vary among bell-like, skewed bell, flattened

bell, or any combination of bell-like forms. This observation is intuitively appealing and a few illustrations will show that it is indeed logical. The histogram can be expected to contain local maxima and minima.



Figure 4

The ideas presented in this section will be used as keystones in those portions of the procedure concerned with analysis of histograms.

D. THE ANALYSIS OF PROJECTED POINTS

As stated earlier, the procedure being developed involves the projection of points on a line and the subsequent examination of the distribution of projected points. The dispersion of points on the line is analyzed through the use of a histogram. Divisions are made between maxima and each segment again analyzed through the use of a new line. Ultimately, adjacent groups are combined and analyzed to correct the possible error that too many divisions were made.

Assuming that the line which maximizes the variance of the points has been found using the non-linear programming method shown previously, the next step is to form a histogram. Fundamental to its construction is a decision concerning the interval length between the equidistant parallel $(n-1)$ -dimensional hyperplanes. It has been determined that, in most cases, dividing the line length spanned

by the two end points into $m/2$ equal intervals, where m is the number of points under consideration, appears to be a reasonable starting interval length. Then if this does not suffice, the interval length yielding more well-defined maxima can be determined through experimentation if desired.

Assuming the efforts in histogram construction are successful, its analysis can now be undertaken. If the plot is unimodal this phase of the procedure is terminated and it can be concluded, for the moment, that all those points comprise one group. As a matter of definition, a plot is considered to be unimodal if there do not exist two separate peaks with a point between them less than eight-tenths (.8) the height of the lesser peak, having zero slope and positive second derivative. Eight-tenths is chosen through the following reasoning: it is advantageous to divide the space into a maximum number of regions to gain information concerning the spatial grouping of points. As the number of groups analyzed simultaneously decreases, the accuracy and informational yield of the analysis increases. The fact that groups will later be combined alleviates any fear of error emanating from splitting conglomerates of points which are actually members of the same fraternity. Thus, it is beneficial to the discrimination process to impose and exercise a rather ruthless division rule, that is, one which dictates divisions readily. On the other hand, an expected result of the non-uniformity of point density is that the histogram of a group contains slight dips and rises

superimposed on an overall bell-like form. Therefore, the ruthlessness in division must be reduced to account for a jagged histogram form. Thus, the cut-off is arbitrarily set to be a dip of twenty percent (20%) the height of the lesser peak. The following illustrations may be helpful in clarifying this concept:

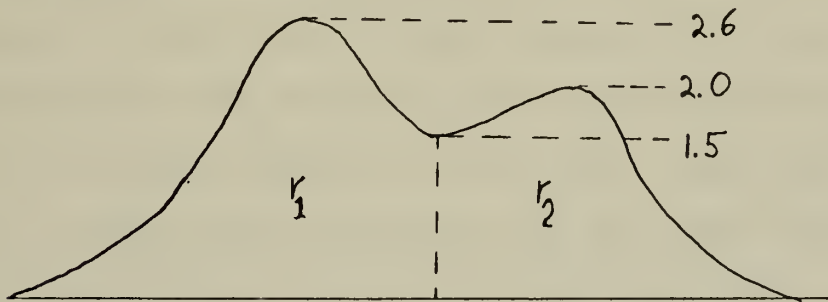


Figure 5

In the above situation, the value of the zero slope point between the two local maxima is sufficiently low to allow a partition of the space into two regions. The second curve permits a three-way partition of n-dimensional space.

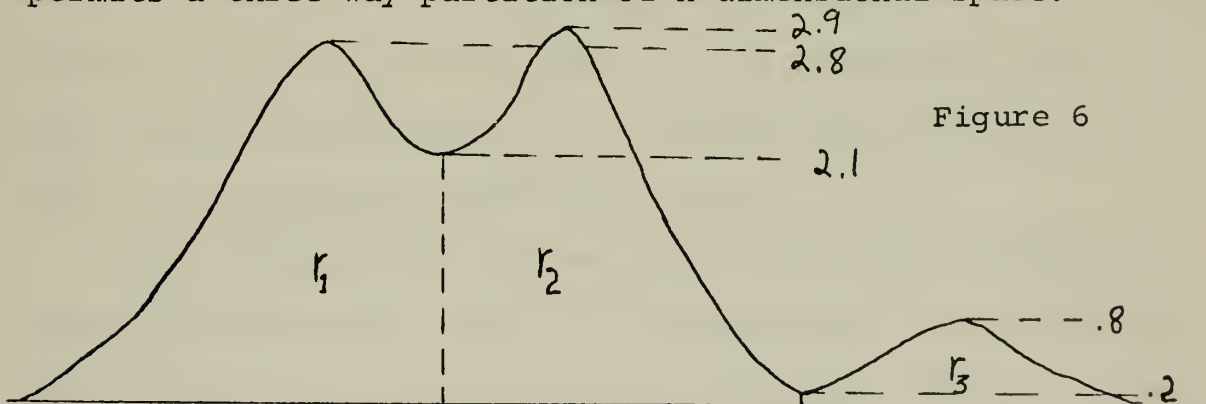


Figure 6

Only one division can be made in the final diagram.

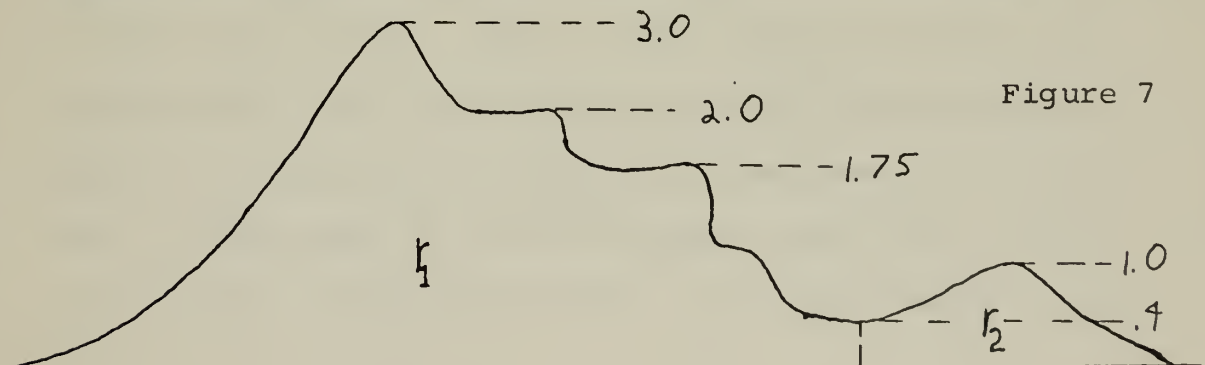


Figure 7

After determining the number of maxima or groups the histogram portrays, the points between any two adjacent peaks must be separated between them. The decision rule for this is to make a division between the maxima at the point c_i where the curve-fitted histogram takes on zero slope and has a positive second derivative as indicated in Figures 5 through 7. Then all of those points to the left of c_i are allocated to the maximum on the left and those to the right of c_i to the maximum on the right. Note that this allocation is made with knowledge that for those cases in which the value of the histogram at the zero slope, positive second derivative point deviates from zero, i.e., for overlapping projections of groups on the line, the method will cause some points to be allocated to a group with which they do not cluster. This inaccuracy is accepted for the present because of the possibility that in the subsequent analysis, those points allocated to the inappropriate group will be split off and later correctly combined.

Then, assuming all of the points have been allocated among several maxima, that is, n -dimensional space has been partitioned into several mutually exclusive and collectively exhaustive regions, the process is continued for each of the regions separately. That is, the entire procedure for finding the best line and analyzing the distribution of points is iterated for each of the regions individually. This is continued for each subregion until none of the regions will subdivide further, or in other words, until

there does not exist a histogram of any of the regions which contains more than one maximum. A summary of this divisive procedure for a typical distribution of points could be represented schematically as follows:

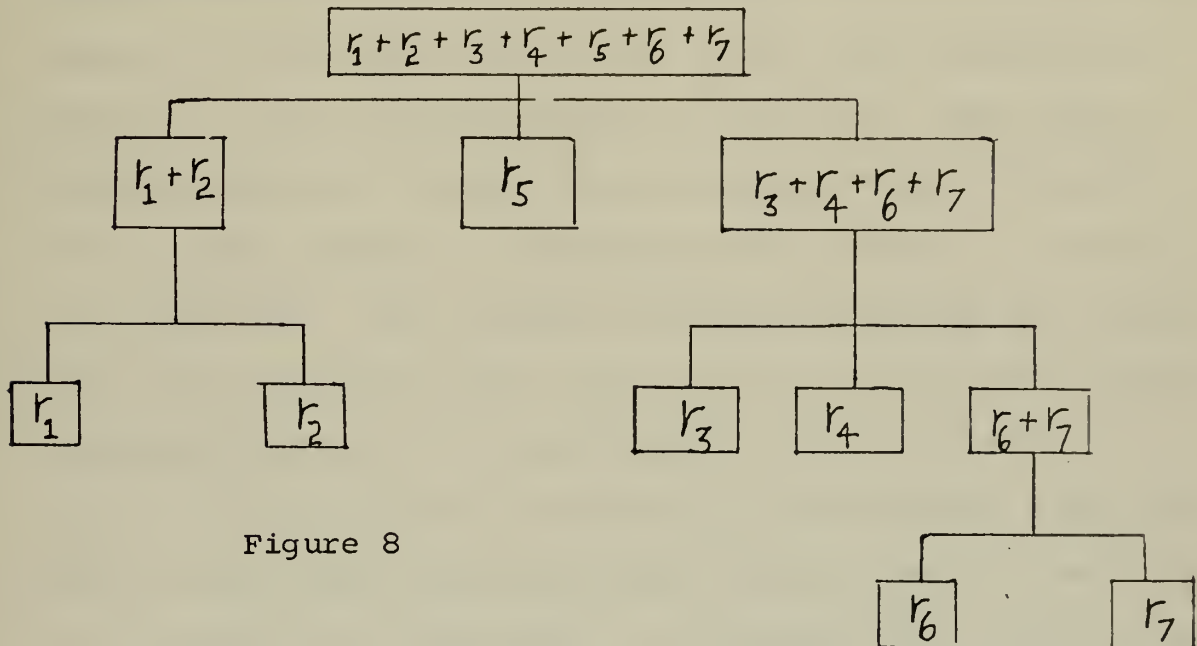


Figure 8

In Figure 8 the procedure has divided n-dimensional space into seven regions. The original maximum variance line contained three distinct maxima and n-dimensional space was partitioned into three regions, $r_1 + r_2$, r_5 , and $r_3 + r_4 + r_6 + r_7$. The maximum variance line for $r_1 + r_2$ and $r_6 + r_7$ showed two maxima and that for $r_3 + r_4 + r_6 + r_7$ contained three. Regions r_1 , r_2 , r_3 , r_4 , r_5 , r_6 , r_7 were unable to further divide because each had unimodal histograms of their individual dispersion of points when projected on the maximum variance line. The result is that n-dimensional space has been partitioned into seven mutually exclusive and collectively exhaustive regions.

E. CORRECTING THE ERROR

As noted previously, in the process of separating points between peaks based solely on the zero slope, positive second derivative criteria, points may have been split from group r_i with which they truly cluster and included in group r_j . If this phenomenon has occurred, there then exist two possible routes these points may have followed:

(1) in later analysis of histograms for group r_j or any of its subgroups, the foreign points were observed to cluster and form a separate peak; they were then split away from the remaining points of r_j and formed a separate group or
(2) due to insufficient clustering, no new maxima containing these points were formed and they were not split from group r_j or any of its subgroups. A procedure for resolving the problem generated by the second alternative has not been developed and must be absorbed as risk. However, the probability of taking the first route is much less than that of taking the second, a result of the aggressive dividing rule.

The implication of the first alternative is that n -dimensional space has been partitioned into too many regions. By combining some of the regions this error can be rectified; that is, adjacent groups must be combined to reduce the number of regions. For definitional purposes those groups adjacent to group r_i will be (1) that group whose centroid is the least number of group r_i 's standard deviations from group r_i 's centroid and (2) that group

containing the point which is closest to any point of group r_i 's, or in other words, that group whose hypersurface is closest to group r_i 's hypersurface. The standard deviation in n-dimensions is calculated solely on the basis of distance from the mean and is thus not a directional quantity. The sample variance is

$$s^2 = \frac{1}{m-1} \sum_{i=1}^m (c_i - \bar{c})^2$$

and the standard deviation is $\sqrt{s^2}$,

where:

m = the number of points

c_i = the coordinates of the i th point = $(c_{i1}, c_{i2}, \dots, c_{in})$

n = the number of dimensions

\bar{c} = the coordinates of the mean point = $(\bar{c}_1, \bar{c}_2, \dots, \bar{c}_n)$

$$= \frac{1}{m} \left(\sum_{i=1}^m c_{i1}, \sum_{i=1}^m c_{i2}, \dots, \sum_{i=1}^m c_{in} \right)$$

Thus the distance of any point, c_i , from the mean is

$$d_i = \left[\left(c_{i1} - \frac{1}{m} \sum_{i=1}^m c_{i1} \right)^2 + \left(c_{i2} - \frac{1}{m} \sum_{j=1}^m c_{j2} \right)^2 + \dots + \left(c_{in} - \frac{1}{m} \sum_{i=1}^m c_{in} \right)^2 \right]^{\frac{1}{2}}$$

Both definitions of adjacent groups are applied in Figure 9.

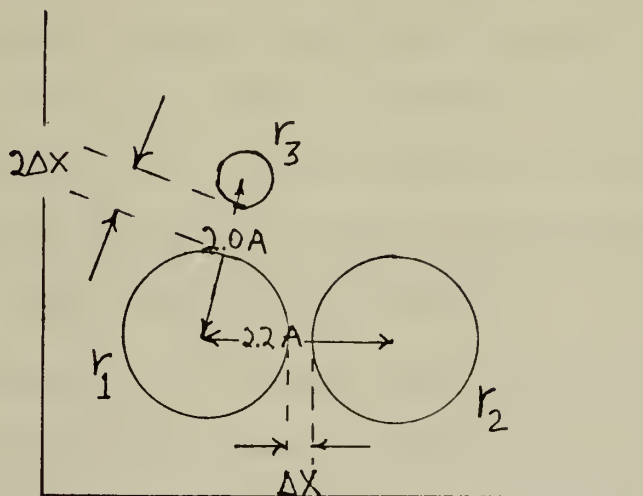


Figure 9

In Figure 9 group r_3 's centroid is only two of group r_1 's standard deviations from group r_1 's centroid and is thus adjacent. Group r_2 is also adjacent to r_1 since its surface is closest to r_1 's.

Thus, for group r_1 with each of its adjacent groups individually in n -dimensional space, the optimal line is determined and the histogram analyzed. The rule for making divisions becomes much more stringent than in the earlier phases of analysis in that partitions of space are not made as readily. The combination of points is separated only if there exists a point between two maxima such that the slope of the curve-fitted histogram is zero, and the value of the curve also equals zero. The reason for this change in policy is that to err by making too few divisions, which appears to be unlikely, is seemingly more acceptable than to err by making too many divisions. This phase is continued until none of the groups will combine with any of its

neighboring groups to form one clump. A typical combination process of the earlier illustrated partitioning of n -dimensional space into seven groups is represented by Figure 10. Initially there are seven groups, r_1, r_2, \dots, r_7 . The symbols in each box represent the group and those immediately below the box represent the adjacent groups of the group designated in the box.

First r_1 is taken with each of its adjacent groups and it is found that it combines with r_2 but not r_7 . In the second step it is determined that $r_1 + r_2$ will not combine with its adjacent group r_4 . When r_3 is examined with each of its adjacent groups it is found that it does not unite with r_5 but will with $r_1 + r_2$. In the third step it can be seen that group $r_1 + r_2 + r_3$ will not combine with its adjacent group, r_5 . Finally it is found that r_5 combines with r_6 . In the last step, r_5 and r_6 are united with r_7 , resulting in three final groups, $r_1 + r_2 + r_3, r_4$ and $r_5 + r_6 + r_7$.

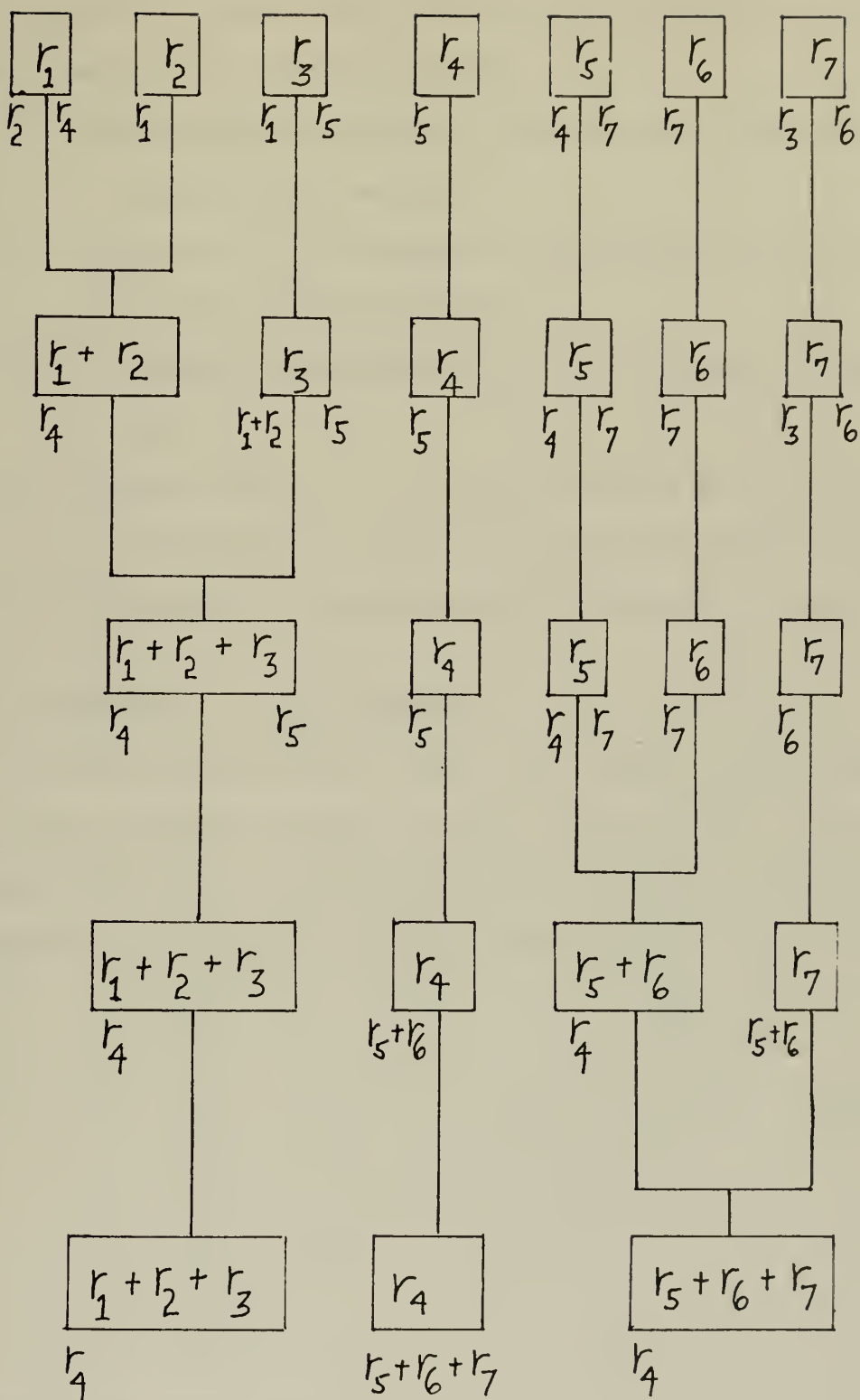


Figure 10

Thus far cluster analysis under this method may be summarized in the following steps.

- (1) Determine the optimal line using the non-linear programming method.
- (2) Construct a histogram representing the distribution of the projected points.
- (3) Analyze the histogram and form a group for each peak.
- (4) Repeat Steps 1, 2, and 3 for each mode and its associated group until none will subdivide.
- (5) Analyze the combinations of adjacent groups.

F. REFINEMENT OF THE PROCEDURE

As mentioned earlier, there are cases in which the line providing maximum variance of the points is not favorable to group detection. One such case involves two side-by-side elongated ellipsoids. Consider Figure 11.

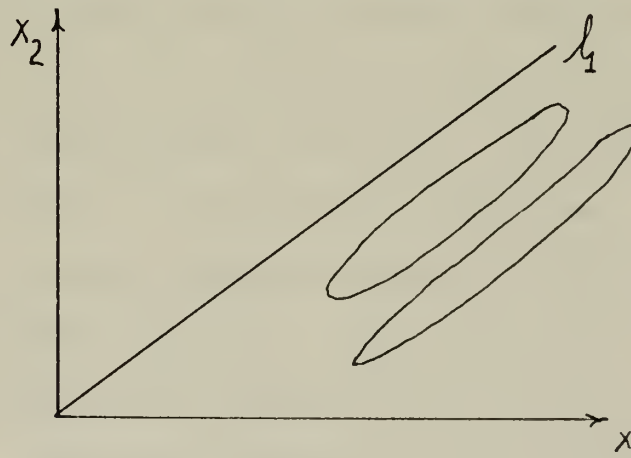


Figure 11

In this instance, the points would be projected on line l_1 and only one group would be detected. Similarly, in n -dimensional space the same problem would be posed, that is,

the maximum variance line would make it possible to detect only one group. It is, therefore, necessary to amend the procedure in a manner which can resolve the elongated ellipsoid case and many of its variations. In most cases it can be seen that there does exist a line onto which the points could be projected which would result in a histogram containing two maxima and hence, the detection of two groups. Figure 12 shows that this is indeed the situation:

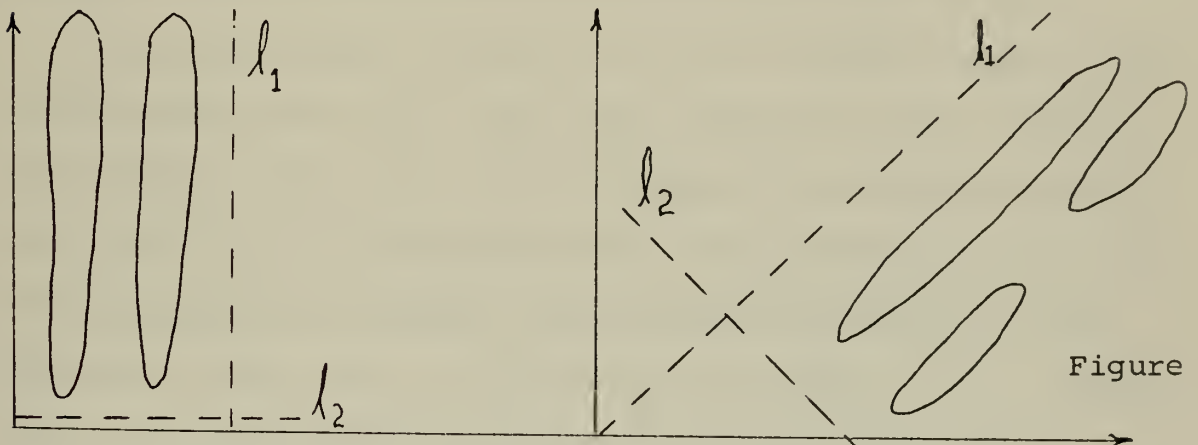


Figure 12

In each diagram, l_1 represents the maximum variance line and l_2 a line which facilitates the detection of two groups. It is interesting to notice that in each case l_2 is perpendicular to l_1 . For the spatial dispersions of points under consideration, that is, for elongated ellipsoids, this result is intuitively appealing.

However, as more characteristics are observed and the space is extended beyond two dimensions, the number of family of lines directionally perpendicular to the maximum variance line, l_1 , increases from one in two space to infinity when the dimension of the space is three or more. In the two dimensional case it can be reasoned that there is only one line family of parallel lines perpendicular to line l_1 .

For the purpose of analyzing point distributions and to be consistent with the non-linear programming solution for the unit vector in the direction of λ_1 , it is reasonable to specify that all lines onto which points are to be projected must pass through the origin. In essence, it is the slope of the line, not its axes intercepts, which is of importance. Similarly, in n-dimensional space the direction or spatial orientation of the line is the only concern.

In the following argument it can be assumed without loss of generality that all lines under discussion pass through the origin. Also λ_1 will designate the maximum variance line and λ_2 will designate the line perpendicular to λ_1 , which successfully handles the elongated ellipsoid. The problem is thus one of ascertaining the line, λ_2 , orthogonal to the maximum variance line which could better discriminate between areas of high point density. It is, therefore, desirable to find a line on which a histogram of the distribution of the projected points will be unimodal if only one group exists and contain two local maxima if two exist. This line, λ_2 , is that line which is perpendicular to λ_1 and on which the variance of the projected points is a relative maximum. That is, considering Figure 12, λ_2 provides a maximum variance between points relative to all other lines perpendicular to the absolute maximum variance line λ_1 , of which there are none. In the three-dimensional analog it can again be reasoned that the line which will facilitate the detection of two groups, if two

exist, is analogous to l_2 . It is again that line which is perpendicular to the absolute maximum variance line, l_1 , and is itself a relative maximum variance line, relative to all those lines perpendicular to l_1 . As the dimension of the space is increased from three to n-dimensions, the line l_2 which is of interest is expected to possess the same characteristics.

The problem of finding such a line can be formulated as a non-linear programming problem of much the same structure as the original method. The only mutation is the addition of a constraint which specifies that the unit vector in the direction of l_2 must be perpendicular to that in the direction of l_1 . Hence, the dot product of the two unit vectors is zero.

The non-linear programming problem is thus:

$$\text{Max } S^2 = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})^2$$

subject to

$$w^t w = 1$$

$$w^t u = 0$$

where: w is a unit vector in the direction of l_2

u is a unit vector in the direction of l_1

(note that all of the coefficients of u are known at this point) and the remaining symbols are unchanged in interpretation from the original programming problem.

The problem now becomes

$$\text{Max } s^2 = \frac{1}{m-1} \sum_{i=1}^m \left[\left(\sum_{j=1}^n c_{ij} w_j \right) - \left(\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n (c_{ij} w_j) \right) \right]^2$$

subject to

$$\begin{aligned} w^t w &= 1 \\ w^t u &= 0 . \end{aligned}$$

The Lagrangian is thus

$$\begin{aligned} F(w_1, w_2, \dots, w_n, \lambda_1, \lambda_2) &= \\ s^2 - \lambda_1 (w^t w - 1) - \lambda_2 (w^t u) &= \\ \frac{1}{m-1} \sum_{i=1}^m \left[\left(\sum_{j=1}^n c_{ij} w_j \right) - \left(\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n c_{ij} w_j \right) \right]^2 \\ - \lambda_1 (w^t w - 1) - \lambda_2 (w^t u) . \end{aligned}$$

As before the partial of the Lagrangian with respect to each variable can be set equal to zero to find the optimal vector w .

Taking the partial with respect to the k^{th} component of w yields

$$0 = \frac{\partial F}{\partial w_k} = \frac{\partial}{\partial w_k} \left[s^2 - \lambda_1 (w^t w - 1) \right] - \frac{\partial}{\partial w_k} (\lambda_2 w^t u),$$

The first term is exactly the same as in the original problem with u replaced by w and thus the equation reduces to

$$0 = \sum_{i=1}^m D_{ik} \sum_{j=1}^n c_{ij} w_j - 2 \lambda_1 w_k - \frac{\partial}{\partial w_k} (\lambda_2 w^t u)$$

$$\text{where } D_{ik} = \frac{2 B_{ik}}{m-1} \text{ and } B_{ik} = c_{ik} - \frac{1}{m} \sum_{i=1}^m c_{ik}.$$

Differentiating the last term with respect to w_k yields only one term, that component of u which is multiplied by w_k .

The derivative of F with respect w_k reduces to

$$(3) \quad 0 = \sum_{i=1}^m D_{ik} \sum_{j=1}^n c_{ij} w_j - \lambda_1 w_k - \lambda_2 u_k, \quad k = 1, \dots, n.$$

There will be n equations of this form, where n is the number of dimensions. Differentiating with respect to the Lagrange multipliers and setting the derivative equal to zero yields the constraints:

$$(4) \quad w^t w = 1$$

$$(5) \quad w^t u = 0.$$

There are now $n+2$ equations in $n+2$ unknowns which can be solved to give the vector w in the direction of ℓ_2 .

The additional step of the procedure is thus to find the line ℓ_2 for each group and observe the resulting histogram of the projected points on ℓ_2 . In this manner the possibility of missing a group is reduced.

The termination of the elongated ellipsoid analysis completes the cluster analysis procedure. In summary the method is described in the following steps.

- (1) Determine the optimal line using the non-linear programming method shown previously.
- (2) Construct a histogram representing the distribution of the projected points.
- (3) Examine the histogram and form a group for each peak.
- (4) Repeat Steps 1, 2, and 3 for each maximum and its associated group until none will sub-divide.
- (5) Analyze the combinations of adjacent groups.
- (6) Examine each group to insure that it cannot be reasonably split into two groups.

It is to be noted that the implementation of the final step need not be confined to the end. For the most favorable result, when the points of a region are projected on their absolute maximum variance line, they should also be projected on their relative maximum variance line. Then, if a division of the region is made, it should be done in accordance with the histogram on that line l_1 or l_2 which most clearly divides the points.

G. INTERPRETATIONS

In light of the earlier discussion of bell-shaped histograms, some interesting observations concerning the various groups can now be made. As noted previously, the more closely the distribution of points in a cluster

approximates a hyperball, the more well defined and bell-like is the curve-fitted histogram representing the distribution of points when orthogonally projected on the optimal line. Thus, simply by observing the histogram for each group, it is possible to make comparisons between groups and to derive useful information about the probable shape of the hypervolume spanned by the points.

Probably the most influential factor in making a judgement will revolve around some within group variance study. In this case, the magnitude of the variance would not be a sufficient comparable statistic between groups; it could easily occur that group r_i and group r_j are equally dense, where density is points per unit volume, but group r_j is larger and contains twice as many points as group r_i . This results in group r_j having a characteristically larger variance, even though both groups may span a similarly shaped hypervolume.

One meaningful comparable statistic could be formed as follows: (1) calculate the standardized variance, i.e., the variance assuming the entire line segment for each group has length one and (2) form a ratio between the height of the mode to that variance, assuming each point in an interval contributes a unit of height to the curve. Observe that the smaller this ratio, the more elongated is the ellipse; in fact, in the limit as the ratio tends to zero, the point distribution approaches a straight line segment.

As a further observation it can be seen that the shape of the ellipsoid yields information concerning the original n specimen variates. Knowledge that a group is an elongated ellipsoid in conjunction with knowledge of its orientation in space is useful in eliciting information about the measured quantities. The process of ascertaining the dominant axis direction, or ellipsoid orientation, is accomplished in the following manner: as before, that line onto which the variance of the orthogonally projected points is a maximum can be determined. This will be that line which is approximately parallel to the longest axis of the ellipsoid and hence, the orientation is specified. Then two particularly interesting cases can occur: (1) the direction of the line may coincide closely with a coordinate axis, that is, the "cigar-shaped" distribution of points may have its longest axis approximately parallel to one of the axis of the coordinate system or (2) the distribution of points may be oriented such that its longest axis lies directionally near the forty-five degree (45°) line between two coordinate axes.

Consider Figure 13, an illustration of the first case.

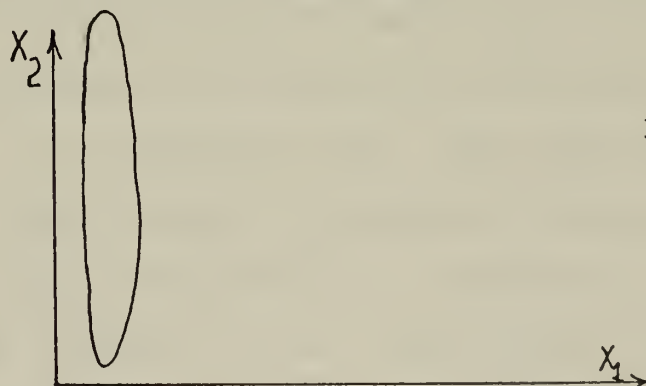


Figure 13

It should be realized that an elongated ellipsoid parallel to an axis indicates that the observations of the sundry characteristics are not varying comparably. That is, one dimension is exerting a relatively powerful force in the determination of point dispersion. Said differently, the observations of one characteristic are varying considerably while observations of the other fluctuate only slightly. Furthermore, an increase in one variate does not affect the other. Thus, for the n-dimensional case, the characteristic represented by the coordinate axis which is parallel to the major axis of the ellipsoid is independent of all other coordinate axes.

The second case, in which a different orientation of the ellipsoid is considered, will first be examined in two-dimensional space. Figure 14 will help to illuminate the discussion.



Figure 14

The elongated ellipsoid has its major axis running in a direction near parallel to the forty-five degree line between the x_1 and x_2 coordinate axes. The implication in this event is that there exists an interaction between the two measured specimen characteristics represented by x_1 and x_2 , hence, the variates are not independent. As an example, suppose that blood pressure and age of humans were being

measured. Generally speaking, blood pressure rises with age and it can be expected that an increase in age would effect an increase in blood pressure. Thus, the expected distribution of points would be in the form of an elongated ellipse whose major axis is positively sloping and directionally close to the forty-five degree line between the representative coordinate axes. An elongated ellipse negatively sloping approximately 135° between two axes indicates a complimentary trade-off between specimen variates; an increase in one variate is characteristically accompanied by a decrease in the other. Extending this concept to n-dimensional space, it can be seen that if the orthogonal projection of a group in the plane formed by any two dimensions results in a "cigar-shaped" distribution of points oriented directionally close to the forty-five degree line between the axes, then the two specimen variates represented by those axes are probably directly related. Said differently, if the unit vector in the direction of the line which maximizes the variance of the projected points of an elongated ellipsoid, as found using the non-linear programming method, contains two direction cosines which are approximately equal, then there may exist a dependence between the underlying variates. In addition, if it is desired to know more precisely the relation between two variates a linear regression analysis can be performed.

H. CAUTION

Although the procedure presented in this paper is capable of handling a wide variety of point dispersion patterns, there still exists many cases in which the techniques of this method are inadequate. For example, consider Figures 15a and 15b.

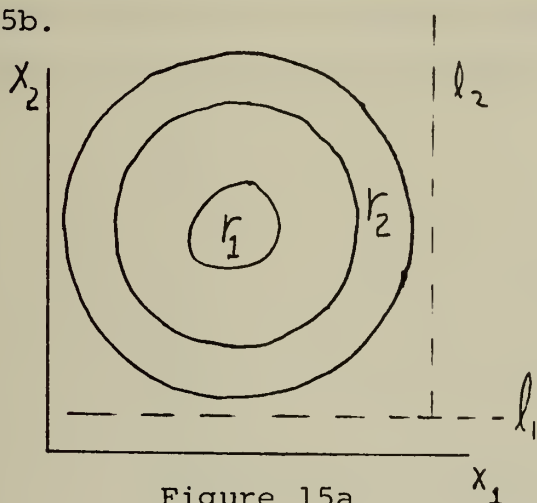


Figure 15a

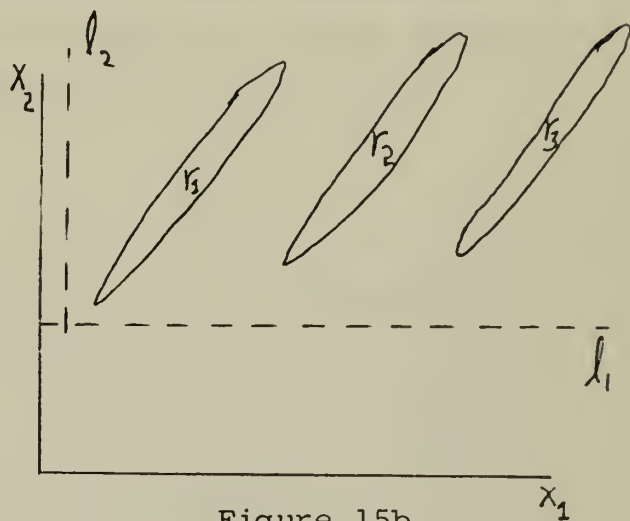


Figure 15b

In Figure 15a it is not possible to detect the smaller group completely surrounded by a "donut-shaped" group. The n -dimensional analog of this is a hyperball (solid sphere) contained in a hypersphere (hollow sphere). Again it can be seen that only one group would be detected. In Figure 15b l_1 is the maximum variance line and l_2 is perpendicular to l_1 . It is possible that there exists precisely that amount of overlap between the group projections on l_1 to cause the histogram to be unimodal. In that case only one group is detected using l_1 . When line l_2 is examined again only one group is detected. In either case A or B it is clear that it is not desirable to term the entire dispersion of points one cluster.

Many of the rules set forth in this work are designed to handle a general case. The user may wish to alter some of these to fit his specific data, his motives for using cluster analysis and his desired accuracy. For example, depending upon the specific criteria selected the two-dimensional point dispersion shown in Figure 16 might be termed as one group or as two by the method presented in this paper.

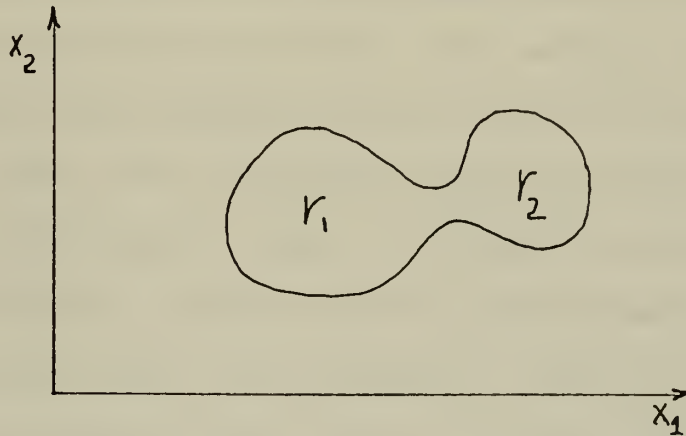


Figure 16

The reader may decide if the point dispersion is one or two groups by relaxing or tightening the constraint for making divisions when combining adjacent groups.

III. EXAMPLE

A. EXPLANATION OF THE DATA

To illustrate the operation of the analysis procedure, an example was chosen in which thirteen variates were measured on 254 specimens. The specimens are naval officers who attended the Naval Postgraduate School, Monterey, during the years 1964 through 1966. The variates are the results of several measurements derived from the Allport, Vernon, and Lindzey (AVL) and the Edwards Personal Preference Schedule (EPPS) psychological tests. Through the theory of psychology and personality, as set forth in Edward Spranger's Types of Men, the two above tests measure a total of twenty-one aspects of human nature, six in the first test and fifteen in the second. To avoid repetition, only thirteen of these variates are used in the cluster analysis example, six from AVL and seven from the EPPS.

The six scores from the AVL reflect the following aspects of human make-up: (1) theoretical, which pertains to the desire to pursue truth, (2) economic, or the interest in what is useful, (3) aesthetic, which involves the interest in form and harmony, (4) social, or the high esteem for love of people and unselfishness, (5) political, or lust for power and (6) religious, which involves the desire for unity or comprehension of the whole.

Those variates chosen from the EPPS measure: (1) achievement, or the ability to accomplish tasks requiring skill and

effort, (2) deference, which involves the acceptance or acknowledgement of the leadership of others, (3) affiliation, or loyalty to friends or groups, (4) intraception, or the ability to analyze the behavior of others, (5) dominance, or the facility to be a leader, (6) nurturance, which is giving of oneself to others who need kindness or understanding and (7) endurance, or the ability to see a job through to completion.

It is desirable to know if the 254 naval officers cluster as one high point density area or as several when the above thirteen variates are measured on each. If there do exist several groups, the characteristics of each could be analyzed and perhaps yield valuable information to the Navy. For example, if most of the officers in one group were those who voluntarily resigned, the AVL and the EPPS could then be used as a predictor to aid the Navy in determining beforehand those officers who would probably not make the service a career. Such a result would indeed have great usefulness within the Navy. On the other hand, if points representing the officers cluster into only one group, the indication is that the AVL and EPPS would not be successful as predictors for officer career patterns.

B. PROCEDURE

The first step in the cluster analysis method is to project the points on their maximum variance line. The Fortran IV program included at the end of this paper first determines the set of non-linear equations to be solved,

then determines their solution, which is in the form of the unit vector in the direction of the maximum variance line, and ultimately projects all points on that line. The only inputs to the program necessary are the coordinates of all the points and the number of dimensions. Using this program it was found that the initial set of fourteen simultaneous non-linear equations to be solved when considering all 254 points are:

$$AU - 2 \lambda U = 0$$

$$\sum_{i=1}^{13} u_i^2 - 1 = 0$$

where

$$U = (u_1, u_2, \dots, u_{13})$$

and

A is the thirteen by thirteen array of numbers defined on the following page.

A =

33.3	-3.9	-6.7	-1.2	5.9	-12.0	3.8	10.5	-.7	47.2	-6.7	6.1	-15.7
3.9	22.6	2.6	3.1	-1.9	.8	-.2	-2.5	-2.8	3.0	-2.1	-3.2	7.5
-6.7	2.6	31.4	2.7	-4.3	17.4	-8.7	-3.0	-1.2	-6.1	4.3	-2.9	7.6
-1.2	3.1	2.7	47.7	-.2	.6	2.6	-3.2	-6.7	-5.2	1.7	.9	12.9
5.9	-1.9	-4.3	-.2	36.1	-6.7	.1	6.6	9.3	-7.1	-1.8	-13.1	-18.5
-12.0	.8	17.4	.6	-6.7	36.1	5.2	-6.0	-5.3	-10.6	6.4	-6.8	21.5
3.8	-.2	-8.7	2.6	.1	-5.2	49.7	14.2	4.5	-4.6	-7.8	-11.7	5.9
10.5	2.5	3.0	-3.2	6.6	-6.0	14.2	67.3	5.5	-.2	-23.0	-7.1	-40.8
-.7	-2.8	-1.2	-6.7	9.3	-5.3	4.5	5.5	99.9	-64.2	-38.7	38.6	-38.3
7.2	3.0	-6.1	-5.2	-7.1	-10.6	-9.6	-.2	-64.2	165.5	3.5	-55.9	-48.7
-6.7	-2.1	4.3	1.7	-1.8	6.4	-7.8	-23.0	-38.7	3.5	94.2	-42.2	6.1
6.1	-3.2	-2.9	.9	13.1	-6.8	-11.7	-7.1	38.6	55.9	-42.2	101.8	-30.6
-15.7	7.5	7.6	12.9	-18.5	21.5	5.9	-40.8	-38.3	-48.7	6.0	-30.6	154.2

Their solution in terms of u_1 through u_{13} is

$$u = (-.32, .72, -.23, -.06, -.08, -.18, -.10, -.07, -.26, \\ -.26, -.16, -.17, -.23)$$

This vector, u , is the unit vector, u , in the direction of the maximum variance line. The 254 points were then projected on the line whose direction is given by the above unit vector. As can be seen by the first histogram in the appendix, the curve-fitted histogram of the projected points is unimodal. Therefore, the points were projected on the relative maximum variance line, which is perpendicular to the absolute maximum variance line. The curve-fitted histogram as shown in the appendix contains two local maxima and a division can be made in accordance with the division rule requiring the local minima to be no greater than eight-tenths the height of the lesser adjacent peaks. Thus, n -dimensional space has been partitioned into two mutually exclusive and collectively exhaustive regions r_1 , containing 142 points and r_2 , containing the remaining 112 points. The next two curves indicate that the region r_1 would not subdivide further, either when its points were projected on the absolute maximum variance line or the relative maximum variance line. When region r_2 was considered it was found that a division was possible when the points were projected on their absolute maximum variance line. Region r_2 was, therefore, partitioned into r_{2a} and r_{2b} ,

with fifty-two points allocated to the former and sixty to the latter.

As shown in the following diagrams, it was not possible to make any further division of any of the regions. Furthermore, region r_{2a} was found to be adjacent to r_{2b} . The two regions were combined into one, r_2 , since the zero slope, positive second derivative point on the curve-fitted histogram had value approximately eighty percent the height of the lesser peak. Similarly, r_1 and r_2 were combined. According to this method, then, the points are one cluster. It is interesting to note that this result is in agreement with several different statistical techniques applied to this same data. These methods include various forms of regression analysis and a discriminate analysis technique developed at The University of California at Berkeley and presented by J. W. Dixon, (Ref. 2).

C. CONCLUSION

The conclusion which must be reached in accordance with this cluster analysis on the given data is that the AVL and EPPS were not shown to be valid for distinguishing among naval officers with respect to the thirteen criteria mentioned. It should be noted, however, that the data used was very narrow-based; that is, it represented only a small portion of the officers in the Navy, namely a fraction of those attending graduate school. If a more representative sample could be obtained, it is indeed feasible that different results would be obtained.

IV. APPENDIX A :
HISTOGRAM ANALYSIS

This appendix contains ten histograms and their curve-fitted approximations. On each diagram is the unit vector specifying the line onto which the points were projected. In those cases where a division was made the division line is shown.

Following the histograms is a computer program which yields sufficient data to facilitate a histogram analysis. The program first determines the set of non-linear equations to be solved for the unit vector, it secondly determines their solution, and thirdly projects the points onto the line specified by the unit vector. Thus, for each histogram shown, a set of non-linear equations was found and solved for the unit vector and the points were subsequently projected on the specified line by the computer program.

All points projected
on their maximum
variance line.

$$\mu = (-.32, .72, -.23, -.06, -.08, \\ -.18, -.10, -.07, -.26, \\ -.26, -.26, -.16, -.17, \\ -.23)$$

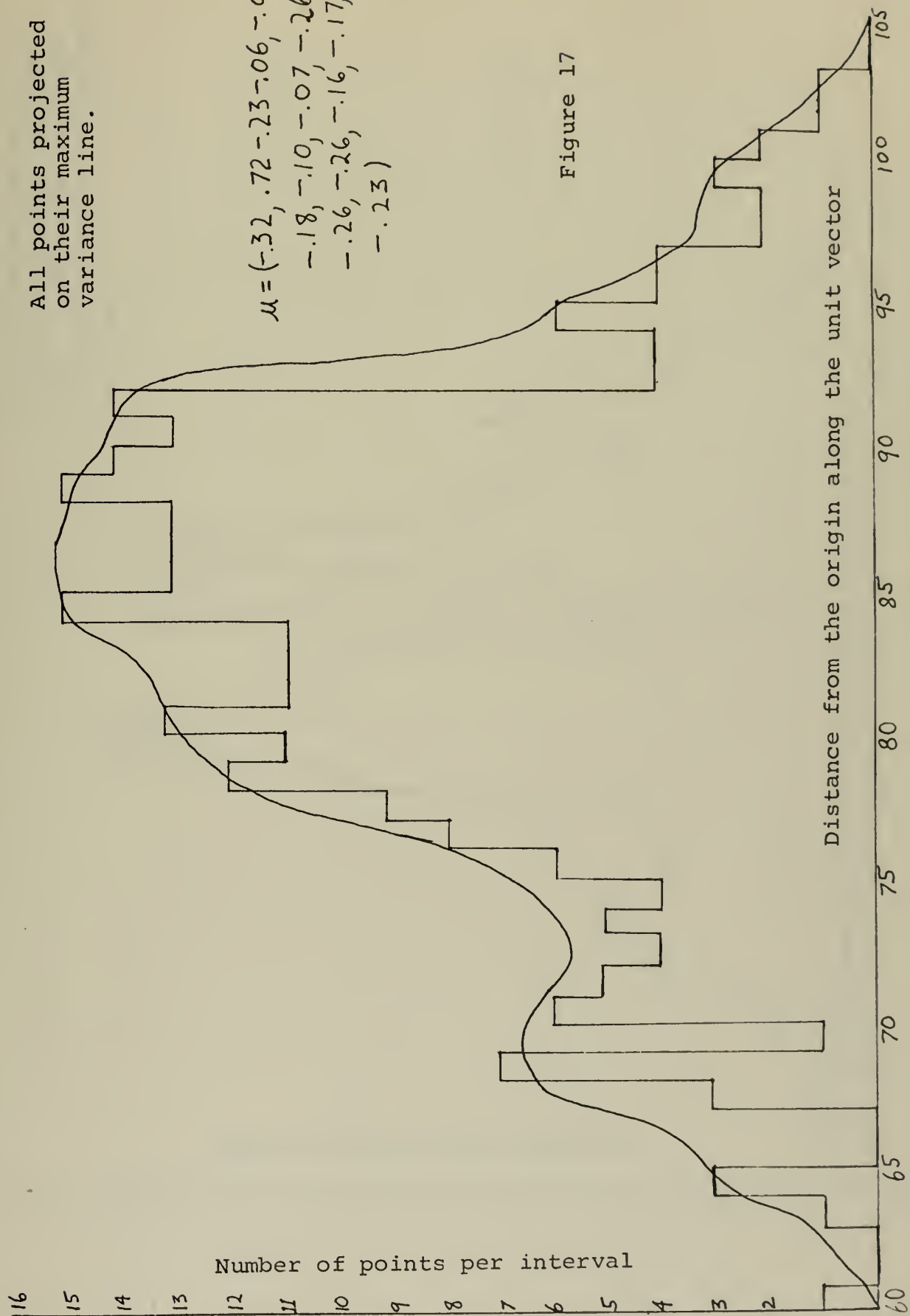
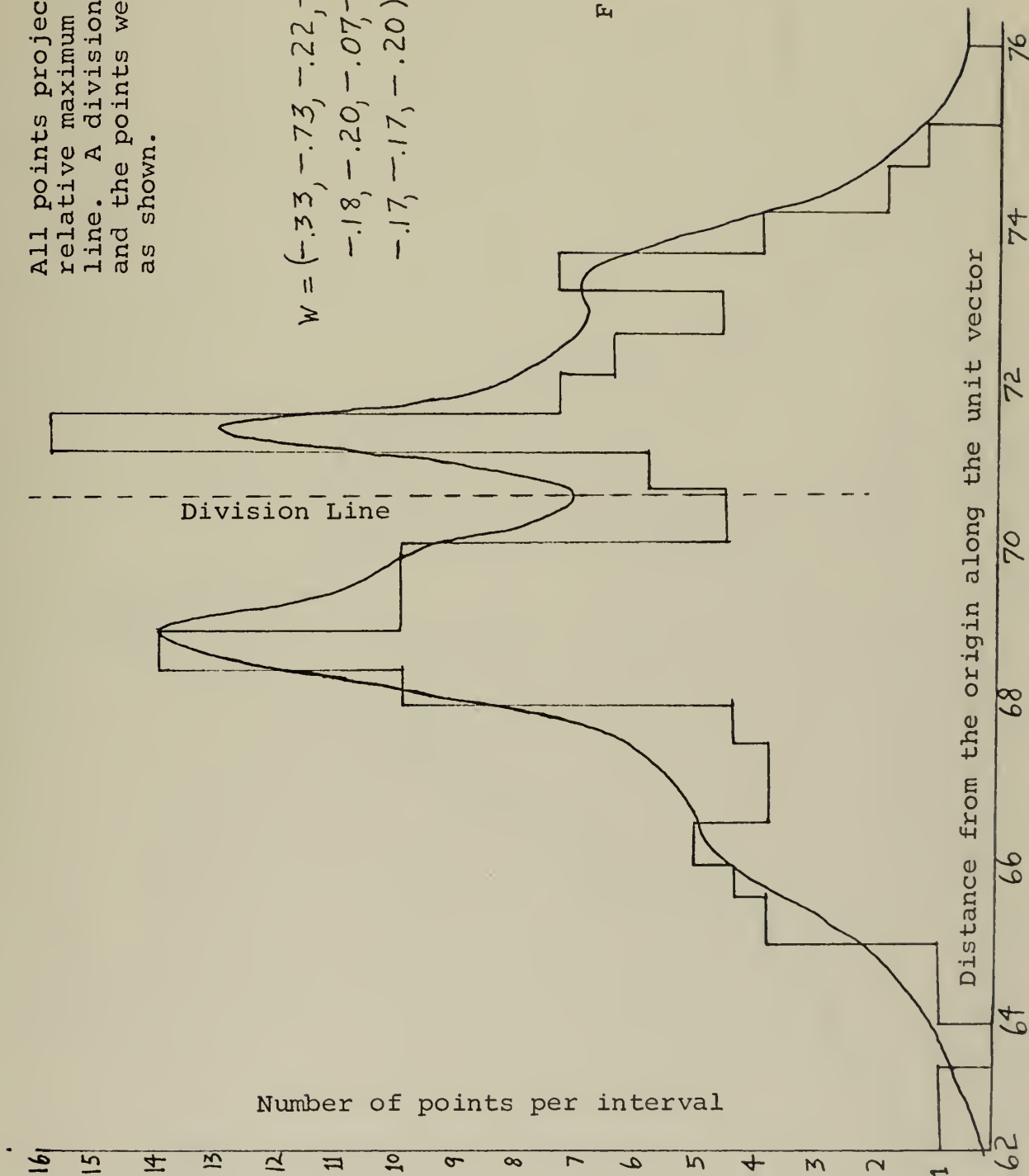


Figure 17

All points projected on the relative maximum variance line. A division was made and the points were separated as shown.

$$W = \begin{pmatrix} -.33, -.73, -.22, -.05, -.06 \\ -.18, -.20, -.07, -.26, -.25 \\ -.17, -.17, -.20 \end{pmatrix}$$

Figure 18



Region r_1 projected on its
maximum variance line.

$$\mu = (.56, -.21, .06, -.33, -.37, \\ .11, .45, .17, .13, .12, \\ -.30, -.01, .01)$$

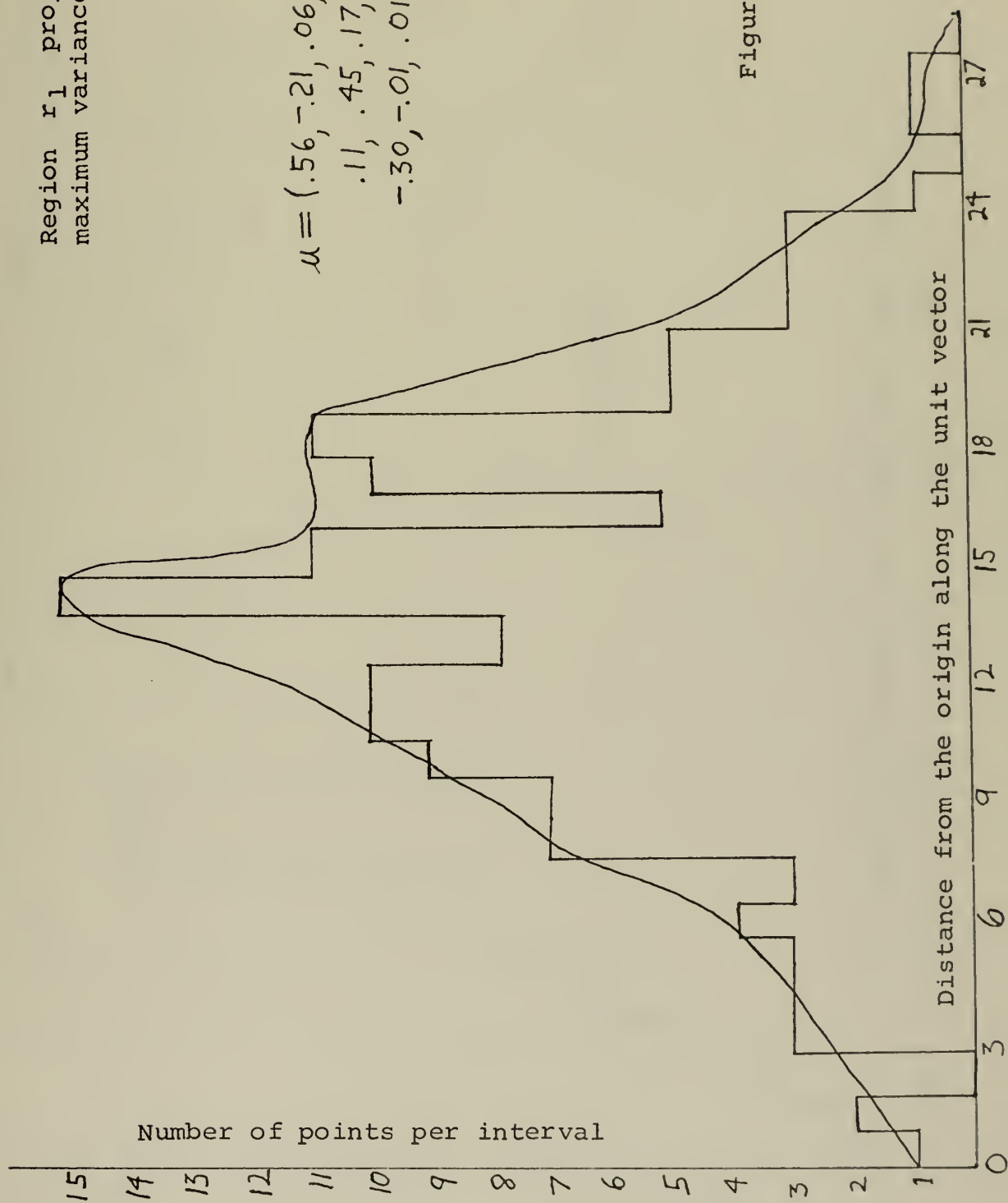
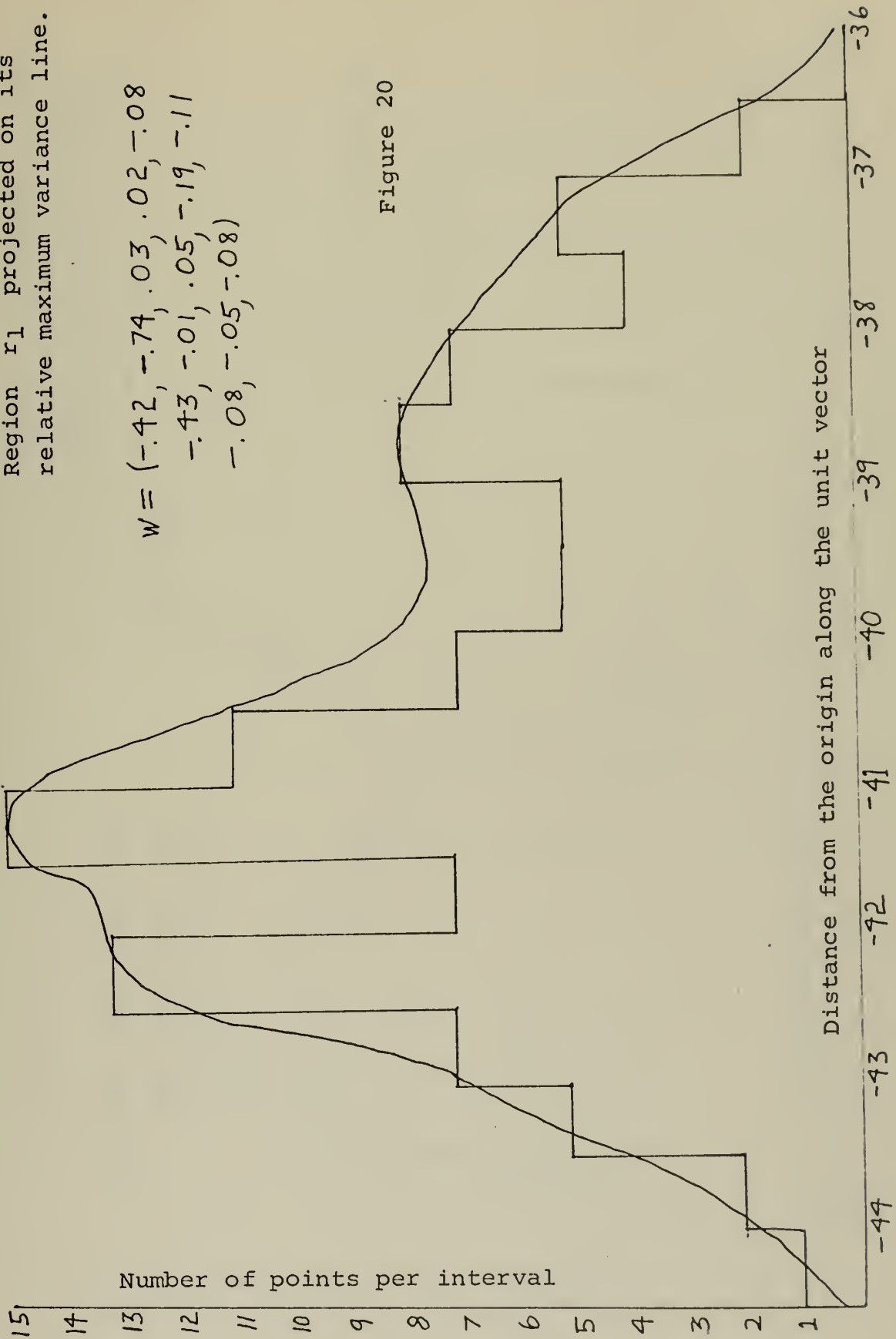


Figure 19

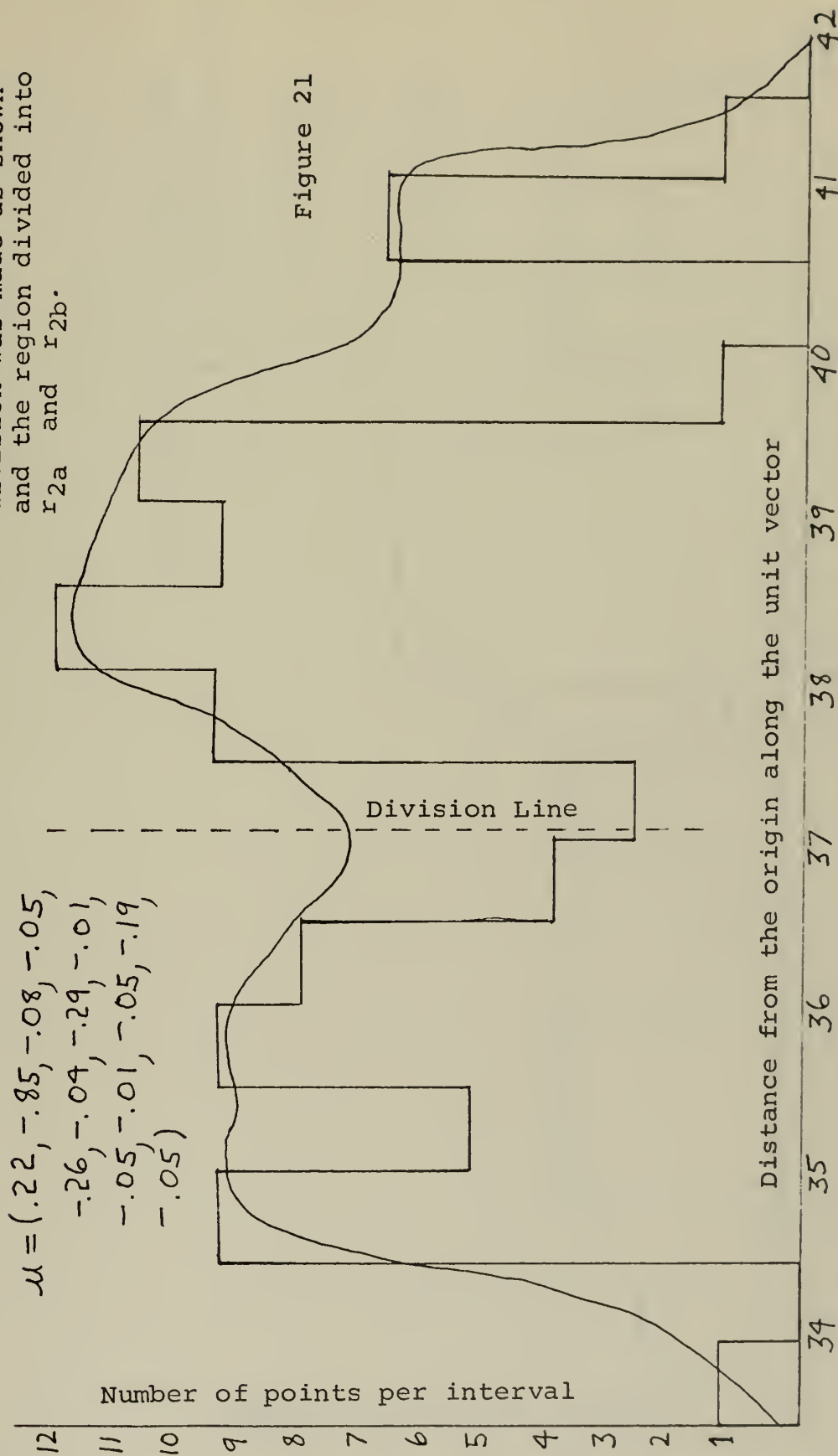
Region r_1 projected on its
relative maximum variance line.

$$W = (-.42, -.74, .03, .02, -.08, \\ -.43, -.01, .05, -.19, -.11, \\ -.08, -.05, -.08)$$

Figure 20



Region r_2 projected on its maximum variance line. A division was made as shown and the region divided into r_{2a} and r_{2b} .



Region r_2 projected on its
relative maximum variance
line.

$$W = (-.25, -.85, -.19, .07, .01, -.02, \\ -.15, .13, -.13, -.13, -.15, -.18, \\ -.17)$$

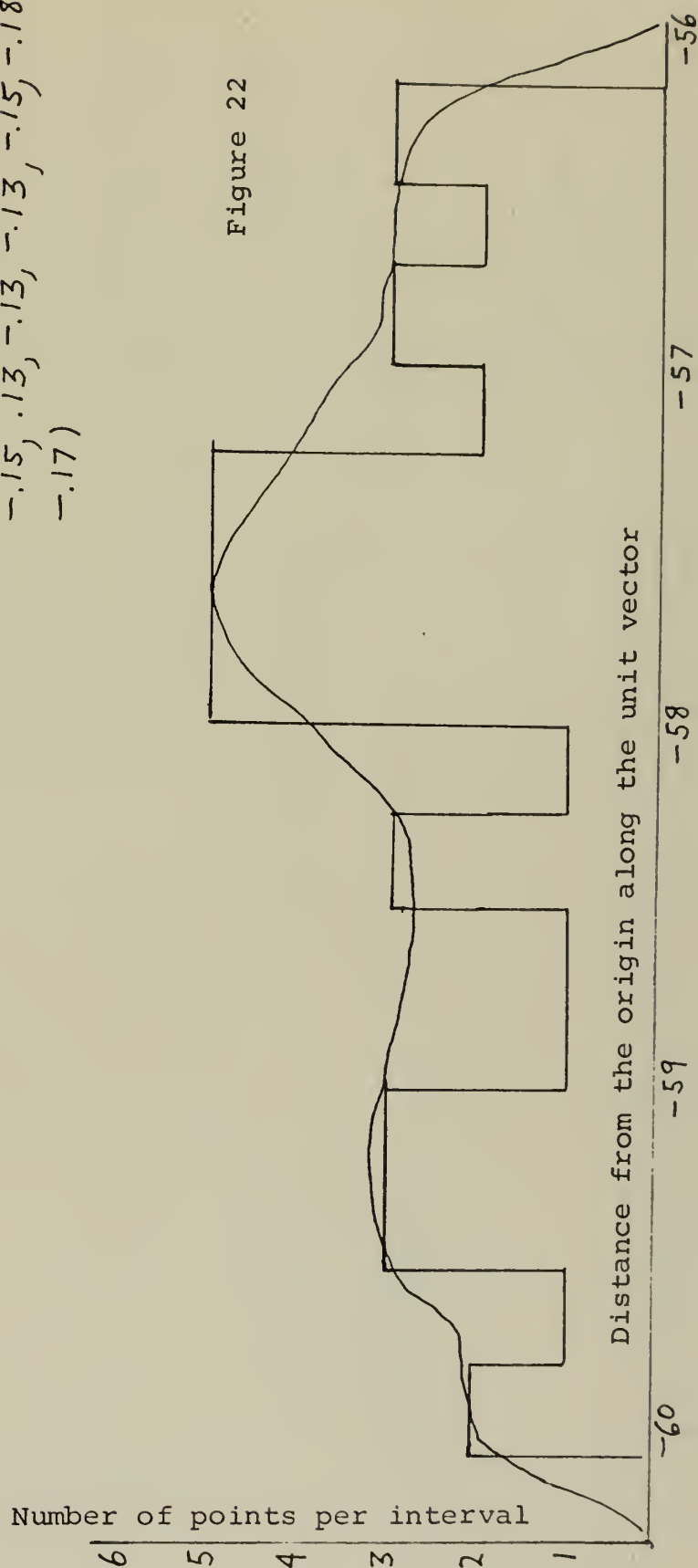


Figure 22

Region r_{2a} projected on its
maximum variance line.

$$\mu = (-.02, -.43, .01, .13, .60, \\ .41, .36, -.32, .03, .02, \\ -.11, -.05, -.03)$$

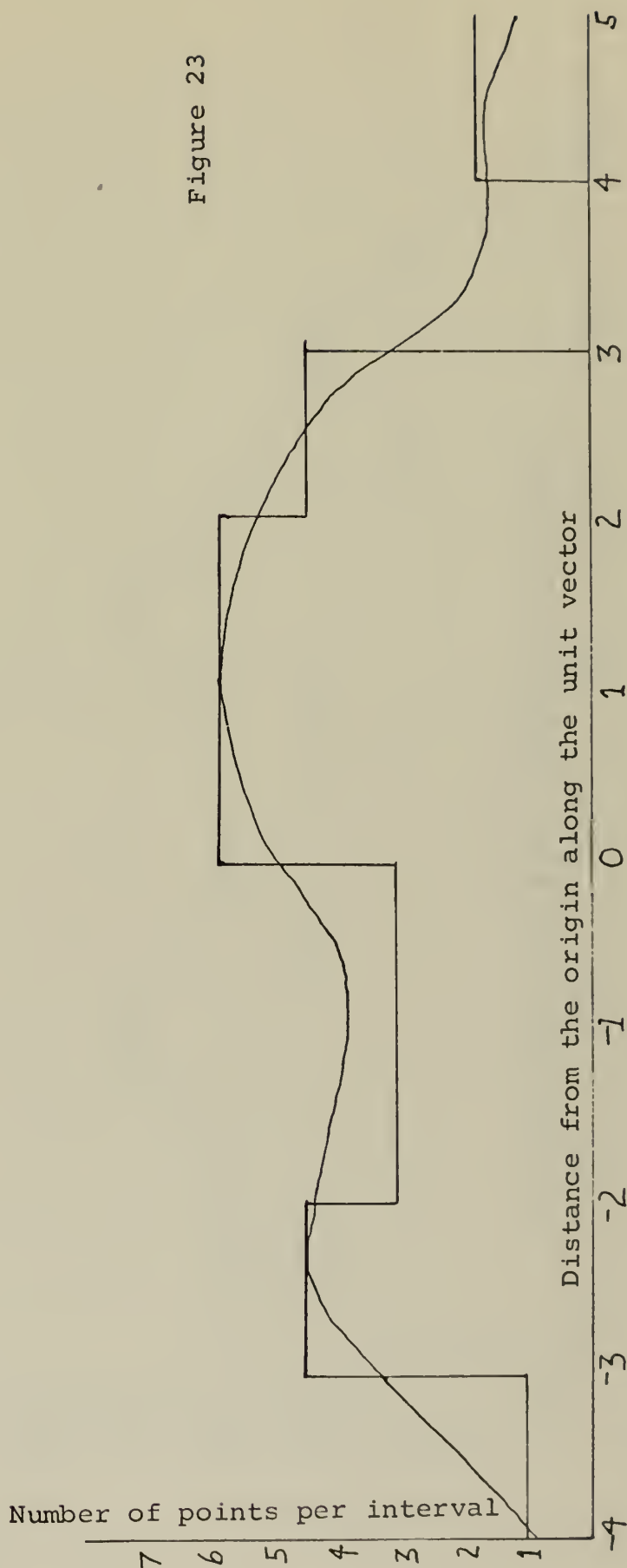


Figure 23

Region r_{2a} projected on its
relative maximum variance line.

$$W = (-.17, -.56, -.31, .24, .35, .30, .10, -.30, -.13, -.16, -.25, -.19, -.18)$$

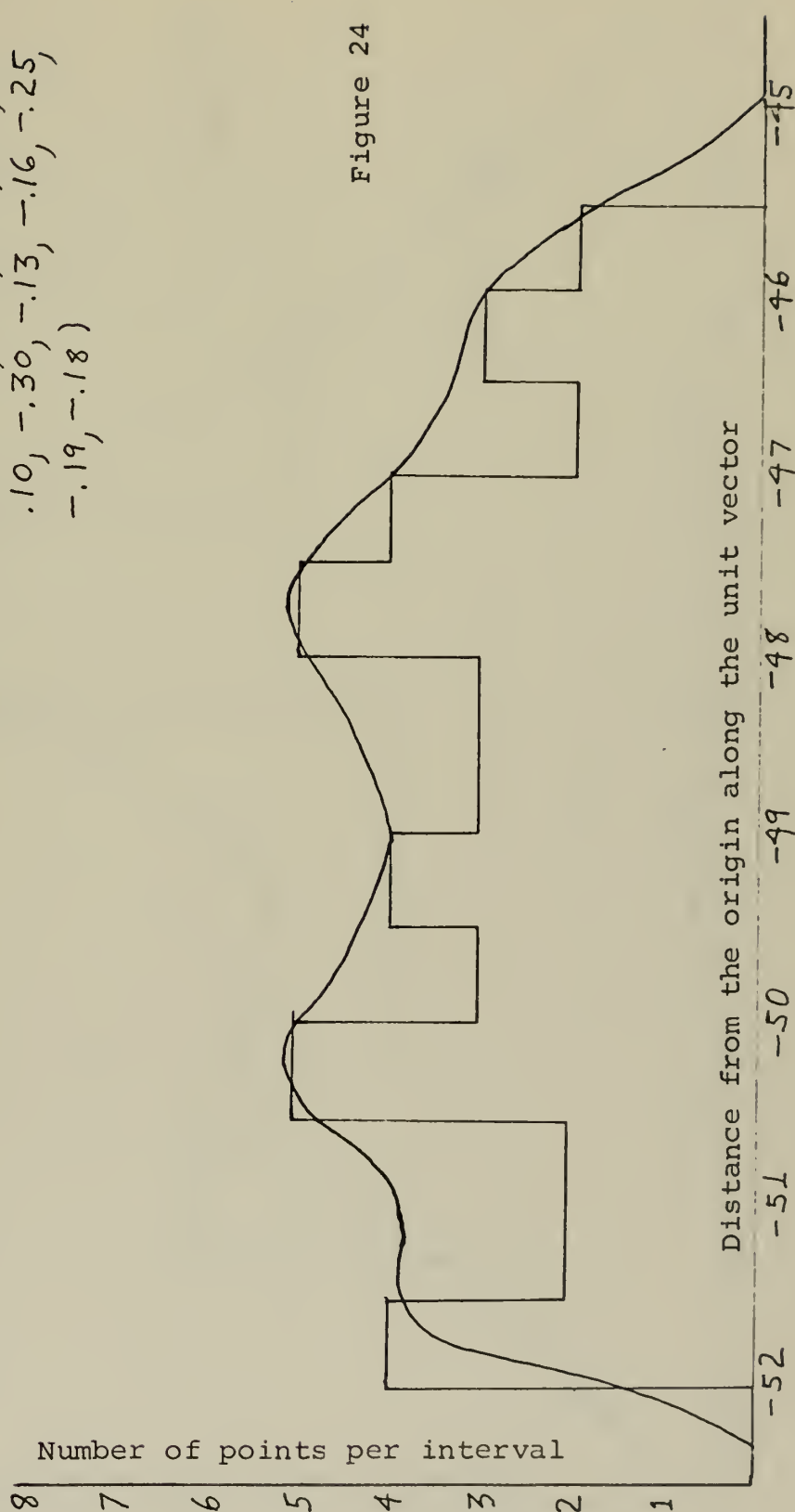
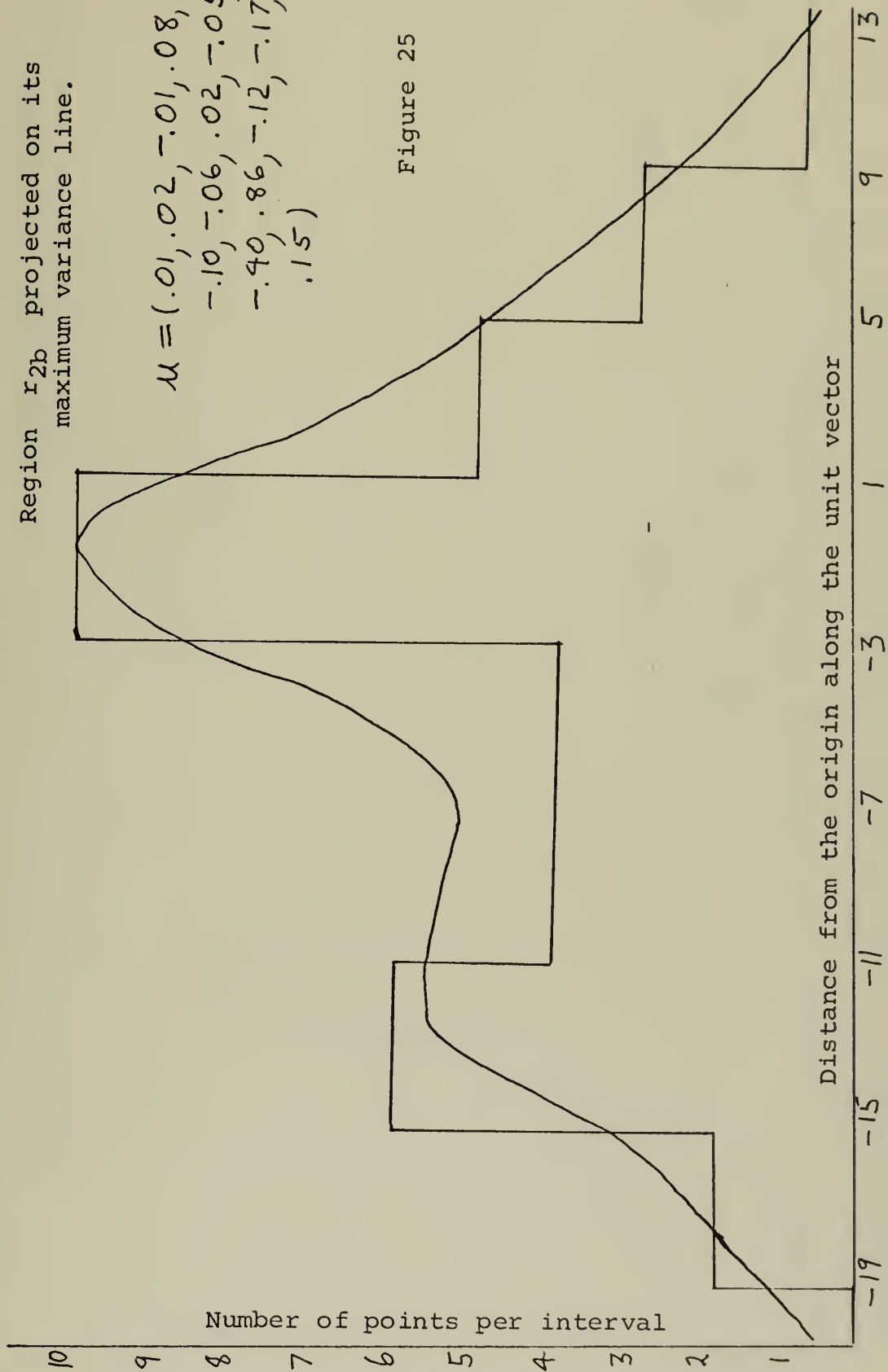


Figure 24

Region r_{2b} projected on its
maximum variance line.

$$\mu = (.01, .02, -.01, .08, \\ -.10, -.06, .02, -.05, \\ -.40, .86, -.12, -.17, \\ .15)$$

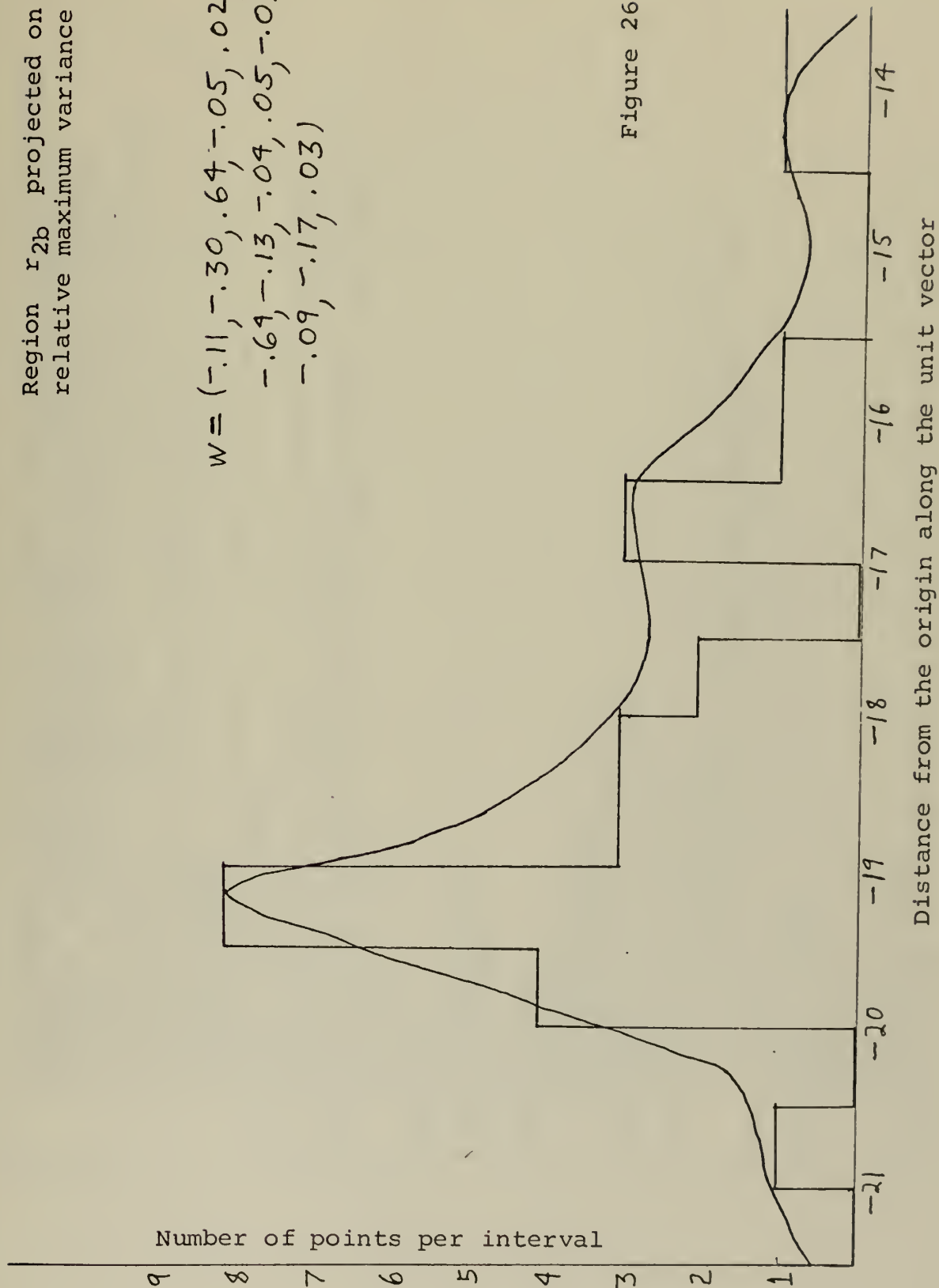
Figure 25



Region r_{2b} projected on its
relative maximum variance line.

$$W = (-.11, -.30, .64, -.05, .02, \\ -.64, -.13, -.04, .05, -.03, \\ -.09, -.17, .03)$$

Figure 26



COMPUTER PROGRAM

```

C THIS PROGRAM IS DESIGNED TO DO SEVERAL PHASES OF THE COMPUTATIONAL
C WORK IN THE CLUSTER ANALYSIS METHOD. IT FIRST DETERMINES THE
C SYSTEM OF N SIMULTANEOUS NON-LINEAR EQUATIONS TO BE SOLVED.
C SECONDLY, IT SOLVES THESE EQUATIONS FOR THE OPTIMAL (MAXIMUM
C VARIANCE, OR RELATIVE MAXIMUM VARIANCE) LINE AND ULTIMATELY
C PROJECTS POINTS IN A UNIT VECTOR IN THE DIRECTION OF THIS LINE.
C DIMENSION UU(13), FF(259)
C DIMENSION X(28)
C DIMENSION C(259,13), B(259,13), U(13,13), D(259,13)
C EXTERNAL EVALUT
C INTEGER POINTER
C COMMON POINTER(30,30), COE(30,30), ISUB(30), U, MM, NN
C READ (5,2001) MM, NN
C FORMAT (2I4)
2001 FM IS THE NUMBER OF POINTS. NN IS THE NUMBER OF DIMENSIONS,
C
2002 READ (5,2002) ((C(I,J), J=1, NN), I=1, MM)
C
C'S ARE THE POINTS. C(I,J), J=1, ..., NN, ARE THE COORDINATES
C OF THE ITH POINT.
C NN=NN+1
K=1
2003 SUM=C. C
C 2010 I=1, MM
2010 SUM=SUM+C(I,K)
C SUM=SUM/MM
J=1
2004 CONTINUE
C U(K,J)=C(I,K)-SUM
C 2020 I=1, MM
2020 B(I,K)=C(I,K)-SUM
C D(I,K)=2.0*B(I,K)/(MM-1.0)
2020 U(K,J)=D(I,K)*C(I,J)+U(K,J)
C J=J+1
2005 IF(J>NN)2004,2004,2005
C K=K+1
C IF(K>NN)2003,2003,2006
C AT THIS POINT ALL OF THE U VALUES HAVE BEEN FOUND. U(K,J)
C IS THE MULTIPLIER OF X(J) IN THE KTH EQUATION.
2006 CONTINUE
2015 WRITE(6,2015) ((U(K,J), J=1, NN), K=1, NN)
C FORMAT (F9.4)
C N=27
C MAXIT=1000
C NUMSIG=4
C IPRINT=1
C 4000 I=1, NN
C X(I)=.26
C X(14)=5.0

```



```

X(28)=5.0
NNNN=NN+2
DO 4050 II=NNNN,N
  X(II)=1.0
  CALL NLNSYS(N,MAXIT,NUMSIG,ISING,IPRINT,EVALUT,X)
  IF (ISING.EQ.0) WRITE (6,2000)
  GO TO 3001
C   AT THIS POINT, THE SET OF EQUATIONS HAS BEEN SOLVED AND POINTS WILL NOW BE
C   AT THIS POINT, THE SET OF EQUATIONS HAS BEEN SOLVED AND POINTS
C   WILL NOW BE PROJECTED ON THE LINE.
2000 FORMAT (1H1,55HSYSTEM CREATED A SINGULARITY NEAR THE FOLLOWING X VA
1LUE)
3001 CONTINUE
SUM=0.0
DO 3022 I=1,NN
  SUM=SUM+X(I)**2
  SUM1=SQRT(SUM)
DO 3006 J=1,NN
  UU(J)=X(J)/SUM1
  WRITE (6,1111)
1111 FORMAT (20H THE UNIT VECTOR IS)
  WRITE (6,1000) (UU(J),J=1,NN)
1000 FORMAT (F8.4)
C   UU(J) IS THE JTH COMPONENT OF THE UNIT VECTOR.
DO 3005 I=1,MM
  SUM=0.0
DO 3010 J=1,NN
  SUM=SUM+C(I,J)*UU(J)
3010 FE(I)=SUM
3005 DO 3023 I=1,MM
  WRITE (6,3024) I
3024 FORMAT (I4)
3015 FE(I)=FE(I)+FE(I)
3023 WRITE (6,3015) FE(I)
C   FE(I) IS A CONSTANT WHICH WHEN MULTIPLIED TIMES THE UNIT VECTOR
C   GIVES THE COORDINATE OF THE ITH POINT WHEN PROJECTED OF THE LINE.
  STOP
END

```

```

CCCCCCCC
SUBROUTINE EVALUT(X,K,F)
IF IT IS DESIRED TO PROJECT POINTS ON THE MAXIMUM VARIANCE LINE,
THE EQUATIONS ARE AS SHOWN BELOW. IT WILL BE NOTED THAT THEIR
FORM IS IDENTICAL TO THAT SHOWN IN THE THESIS BODY. THERE ARE
2NN+1 EQUATIONS, WHERE NN IS THE NUMBER OF DIMENSIONS, THE FIRST
NN OF THESE ARE THE FIRST DERIVATIVES OF THE LAGRANGIAN, THE
SECOND NN ARE THE SECOND DERIVATIVES OF THE LAGRANGIAN AND
THE FINAL EQUATION IS THE CONSTRAINT EQUATION. NOTICE THAT

```



```

CCCCCCCCCCCCC
IT IS SPECIFIED THAT ALL SECOND DERIVATIVES BE LESS THAN ZERO,
IN THE EVENT THAT IT IS DESIRED TO PROJECT POINTS ON THE RELATIVE
MAXIMUM VARIANCE LINE, IT IS FIRST NECESSARY TO FIND THE UNIT
VECTOR IN THE DIRECTION OF THE MAXIMUM VARIANCE LINE, THEN,
A VECTOR MUST BE ADDED TO EACH OF THE FIRST NN EQUATIONS OF THE FORM -
-U(1)*X(2) FOR THE ITH EQUATION, I=1,000,NN, WHERE UU(1)
IS THE ITH COMPONENT OF THE UNIT VECTOR, IN THE DIRECTION OF THE
MAXIMUM VARIANCE LINE. IN ADDITION, ANOTHER EQUATION MUST BE
ADDED TO SATISFY THE CONSTRAINT THAT THE RELATIVE MAXIMUM
VARIANCE THIS EQUATION WILL BE OF THE FORM, SUM OVER I OF
+UU(1)*X(1), I=1,000,NN
DIMENSION U(13,13)
INTEGER POINTER
COMMON POINTER(30,30),COE(30,30),ISUB(30),U,MM,NN
IF(K-1) 1001,1001,1002
F=(U(1,1)-2.0*X(14))*X(1)+U(1,2)*X(2)+U(1,3)*X(3)+U(1,4)*X(4)+U(1,
15)*X(5)+U(1,6)*X(6)+U(1,7)*X(7)+U(1,8)*X(8)+U(1,9)*X(9)+U(1,10)*X(
210)+U(1,11)*X(11)+U(1,12)*X(12)+U(1,13)*X(13)
GO TO 1004
1001 IF(K-2) 1003,1003,1005
F=U(2,1)*X(1)+U(2,2)*X(2)+U(2,3)*X(3)+U(2,4)*X(4)+U(2,
15)*X(5)+U(2,6)*X(6)+U(2,7)*X(7)+U(2,8)*X(8)+U(2,9)*X(9)+U(2,10)*X(
210)+U(2,11)*X(11)+U(2,12)*X(12)+U(2,13)*X(13)
GO TO 1004
1002 IF(K-3) 1006,1006,1007
F=U(3,1)*X(1)+U(3,2)*X(2)+U(3,3)*X(3)+U(3,4)*X(4)+U(3,
15)*X(5)+U(3,6)*X(6)+U(3,7)*X(7)+U(3,8)*X(8)+U(3,9)*X(9)+U(3,10)*X(
210)+U(3,11)*X(11)+U(3,12)*X(12)+U(3,13)*X(13)
GO TO 1004
1003 IF(K-4) 1008,1008,1009
F=U(4,1)*X(1)+U(4,2)*X(2)+U(4,3)*X(3)+U(4,4)*X(4)+U(4,
15)*X(5)+U(4,6)*X(6)+U(4,7)*X(7)+U(4,8)*X(8)+U(4,9)*X(9)+U(4,10)*X(
210)+U(4,11)*X(11)+U(4,12)*X(12)+U(4,13)*X(13)
GO TO 1004
1004 IF(K-5) 1010,1010,1011
F=U(5,1)*X(1)+U(5,2)*X(2)+U(5,3)*X(3)+U(5,4)*X(4)+U(5,5)*X(5)+U(5,
15)*X(5)+U(5,6)*X(6)+U(5,7)*X(7)+U(5,8)*X(8)+U(5,9)*X(9)+U(5,10)*X(
210)+U(5,11)*X(11)+U(5,12)*X(12)+U(5,13)*X(13)
GO TO 1004
1005 IF(K-6) 1012,1012,1013
F=U(6,1)*X(1)+U(6,2)*X(2)+U(6,3)*X(3)+U(6,4)*X(4)+U(6,5)*X(5)+U(6,
15)*X(5)+U(6,6)*X(6)+U(6,7)*X(7)+U(6,8)*X(8)+U(6,9)*X(9)+U(6,10)*X(
210)+U(6,11)*X(11)+U(6,12)*X(12)+U(6,13)*X(13)
GO TO 1004
1006 IF(K-7) 1014,1014,1015
F=U(7,1)*X(1)+U(7,2)*X(2)+U(7,3)*X(3)+U(7,4)*X(4)+U(7,5)*X(5)+U(7,

```



```

16)*X(6)+(U(7,7)-2.0*X(14))*X(7)+U(7,8)*X(8)+U(7,9)*X(9)+U(7,10)*X(
210)+U(7,11)*X(11)+U(7,12)*X(12)+U(7,13)*X(13)
GO TO 1004
1015 IF (K-8) 1016,1016,1017
1016 F=U(8,1)*X(1)+U(8,2)*X(2)+U(8,3)+U(8,4)+X(4)+U(8,5)*X(5)+U(8,
16)*X(6)+U(8,7)*X(7)+(U(8,8)-2.0*X(14))*X(8)+U(8,9)*X(9)+U(8,10)*X(
210)+U(8,11)*X(11)+U(8,12)*X(12)+U(8,13)*X(13)
GO TO 1004
1017 IF (K-9) 1018,1018,1019
1018 F=U(9,1)*X(1)+U(9,2)*X(2)+U(9,3)+X(3)+U(9,4)*X(4)+U(9,5)*X(5)+U(9,
16)*X(6)+U(9,7)*X(7)+U(9,8)*X(8)+U(9,9)-2.0*X(14))*X(9)+U(9,10)*X(
210)+U(9,11)*X(11)+U(9,12)*X(12)+U(9,13)*X(13)
GO TO 1004
1019 IF (K-10) 1020,1020,1021
1020 F=U(10,1)*X(1)+U(10,2)*X(2)+U(10,3)*X(3)+U(10,4)*X(4)+U(10,5)*X(5)
1+U(10,6)*X(6)+U(10,7)*X(7)+U(10,8)*X(8)+U(10,9)*X(9)+U(10,10)-2.0
2*X(14))*X(10)+U(10,11)*X(11)+U(10,12)*X(12)+U(10,13)*X(13)
GO TO 1004
1021 IF (K-11) 1022,1022,1023
1022 F=U(11,1)*X(1)+U(11,2)*X(2)+U(11,3)*X(3)+U(11,4)*X(4)+U(11,5)*X(5)
1+U(11,6)*X(6)+U(11,7)*X(7)+U(11,8)*X(8)+U(11,9)*X(9)+U(11,10)*X(10
2)+U(11,11)-2.0*X(14))*X(11)+U(11,12)*X(12)+U(11,13)*X(13)
GO TO 1004
1023 IF (K-12) 1024,1024,1025
1024 F=U(12,1)*X(1)+U(12,2)*X(2)+U(12,3)*X(3)+U(12,4)*X(4)+U(12,5)*X(5)
1+U(12,6)*X(6)+U(12,7)*X(7)+U(12,8)*X(8)+U(12,9)*X(9)+U(12,10)*X(10
2)+U(12,11)*X(11)+U(12,12)-2.0*X(14))*X(12)+U(12,13)*X(13)
GO TO 1004
1025 IF (K-13) 1026,1026,1027
1026 F=U(13,1)*X(1)+U(13,2)*X(2)+U(13,3)*X(3)+U(13,4)*X(4)+U(13,5)*X(5)
1+U(13,6)*X(6)+U(13,7)*X(7)+U(13,8)*X(8)+U(13,9)*X(9)+U(13,10)*X(10
2)+U(13,11)*X(11)+U(13,12)*X(12)+U(13,13)*X(13)
GO TO 1004
1027 IF (K-14) 1028,1028,1029
1028 F=U(14,1)-2.0*X(14)+X(15)**2
GO TO 1004
1029 IF (K-15) 1030,1030,1031
1030 F=U(15,1)-2.0*X(14)+X(15)**2
GO TO 1004
1031 IF (K-16) 1032,1032,1033
1032 F=U(16,1)-2.0*X(14)+X(17)**2
GO TO 1004
1033 IF (K-17) 1034,1034,1035
1034 F=U(17,1)-2.0*X(14)+X(18)**2
GO TO 1004
1035 IF (K-18) 1036,1036,1037
1036 F=U(18,1)-2.0*X(14)+X(19)**2
GO TO 1004

```



```

1037 IF(K-19) 1038,1038,1039
1038 F=U(6,6)-2.0*X(14)+X(20)**2
      GO TO 1004
1039 IF(K-20) 1040,1040,1041
1040 F=U(7,7)-2.0*X(14)+X(21)**2
      GO TO 1004
1041 IF(K-21) 1042,1042,1043
1042 F=U(8,8)-2.0*X(14)+X(22)**2
      GO TO 1004
1043 IF(K-22) 1044,1044,1045
1044 F=U(9,9)-2.0*X(14)+X(23)**2
      GO TO 1004
1045 IF(K-23) 1046,1046,1047
1046 F=U(10,10)-2.0*X(14)+X(24)**2
      GO TO 1004
1047 IF(K-24) 1048,1048,1049
1048 F=U(11,11)-2.0*X(14)+X(25)**2
      GO TO 1004
1049 IF(K-25) 1050,1050,1051
1050 F=U(12,12)-2.0*X(14)+X(26)**2
      GO TO 1004
1051 IF(K-26) 1052,1052,1053
1052 F=U(13,13)-2.0*X(14)+X(27)**2
      GO TO 1004
1053 F=X(1)**2+X(2)**2+X(3)**2+X(4)**2+X(5)**2+X(6)**2+X(7)**2+X(8)**2+
      X(9)**2+X(10)**2+X(11)**2+X(12)**2+X(13)**2-1.0
1004 RETURN
      END

```

```

CCCCCCCCCCCCCCCCCCCC
SUBROUTINE NLNSYS(N,MAXIT,NUMSIG,ISING,IPRINT,EVALUT,X)
      *****
      SUBROUTINE NLNSYS
      IDENTIFICATION - NLNSYS (SOLUTION OF SIMULTANEOUS NON-LINEAR
                        EQUATIONS).
      ID              - C4-NPG-NLNSYS (FORTRAN IV G).
      CATEGORY       - MATHEMATICAL SUBROUTINE
      PROGRAMMER    - D. D. PURCELL
      DATE          - FEB, 1969
      PURPOSE
      TO SOLVE A SYSTEM OF N NON-LINEAR EQUATIONS. THESE MUST BE
      SUPPLIED BY THE USER IN A SUBROUTINE REFERENCED BY THE MAIN-
      LINE CALLING ROUTINE.
      *****
      NLNS1430
      NLNS0009
      NLNS0010
      NLNS0020
      NLNS0030
      NLNS0040
      NLNS0050
      NLNS0060
      NLNS0070
      NLNS0080
      NLNS0090
      NLNS0100
      NLNS0110
      NLNS0120
      NLNS0130
      NLNS0140
      NLNS0150
      NLNS0160

```


NLNS0650
NLNS0660
NLNS0670
NLNS0680
NLNS0690
NLNS0700
NLNS0710
NLNS0720
NLNS0730
NLNS0740
NLNS0750
NLNS0760
NLNS0770
NLNS0780
NLNS0790
NLNS0800
NLNS0810
NLNS0820
NLNS0830
NLNS0840
NLNS0850
NLNS0860
NLNS0870
NLNS0880
NLNS0890
NLNS0900
NLNS0910
NLNS0920
NLNS0930
NLNS0940
NLNS0950
NLNS0960
NLNS0970
NLNS0980
NLNS0990
NLNS1000
NLNS1010
NLNS1020
NLNS1030
NLNS1040
NLNS1050
NLNS1060
NLNS1070
NLNS1080
NLNS1090
NLNS1100
NLNS1110
NLNS1120

THE PARAMETERS ARE RESPECTIVELY - SOLUTION VECTOR GUESS,
NUMBER OF FUNCTION TO BE RETURNED, AND THE RETURN PAR-
AMETER REPRESENTING THE VALUE OF THIS FUNCTION FOR THIS
APPROXIMATION.

DIMENSION X(2)
IF (K.EQ.1) GO TO 1
GO TO 2

1 F=2.71828183*(.920422528*(EXP(2.*X(1)-1.0)-1.0)+X(2)/

3.14159265-2.*X(1)

GO TO 3

2 F=.5*SIN(X(1)*X(2))-X(2)/12.5663706-X(1)/2.

3 RETURN

END

SPACE REQUIRED
NLNSYS PLUS BAKSUB REQUIRE 2400 BYTES AND THE SUPPORT SYS-
TEMS ROUTINES PLUS BUFFERS REQUIRE ABOUT 43,000 BYTES.

CAUTIONS TO USER

THE STARTING GUESS MUST BE WITHIN THE REGION NECESSARY FOR
CONVERGENCE OR THE PROCESS WILL 'BLOW UP'. THIS REGION HAS
NO EASY DEFINITION AND IT MAY BE NECESSARY TO TRY MORE THAN
ONE STARTING GUESS. ALSO, THE USER SUBROUTINE MUST BE TYPED
'EXTERNAL' IN THE CALLING PROGRAM.

EQUIPMENT CONFIGURATION

IBM 260

SUBROUTINE AND FUNCTION SUBPROGRAMS REQUIRED

EVALUT (USER WRITTEN AND USER NAMED)

BAKSUB (SUPPLIED WITH NLNSYS)

STANDARD SYSTEMS SUBPROGRAMS.

REMARKS

THIS PROCEDURE SOLVES A SYSTEM OF N SIMULTANEOUS NON-
LINEAR EQUATIONS. THE METHOD IS ROUGHLY QUADRATICALLY CON-
VERGENT AND REQUIRES ONLY $((N+2)/2)+(3*N/2)$ FUNCTION
EVALUATIONS PER ITERATIVE STEP AS COMPARED WITH $(N+2)*N$
EVALUATIONS FOR NEWTON'S METHOD. THIS RESULTS IN A SAVINGS
OF COMPUTATIONAL EFFORT FOR SUFFICIENTLY COMPLICATED FUNC-
TIONS. A DETAILED DESCRIPTION OF THE GENERAL METHOD AND
TECHNIQUE OF CONVERGENCE ARE INCLUDED IN (2).² BASICALLY THE
PROCEDURE CONSISTS IN EXPANDING THE FIRST EQUATION IN A TAY-
LOR SERIES ABOUT THE STARTING GUESS, RETAINING ONLY LINEAR
TERMS, EQUATING TO ZERO AND SOLVING FOR ONE VARIABLE, SAY
X(K), AS A LINEAR COMBINATION OF THE REMAINING N-1 VAR-
IABLES. IN THE SECOND EQUATION, X(K) IS ELIMINATED BY RE-
PLACING IT WITH ITS LINEAR REPRESENTATION FOUND ABOVE, AND
AGAIN THE PROCESS OF EXPANDING THROUGH LINEAR TERMS, AND
EQUATING TO ZERO AND SOLVING FOR ONE VARIABLE IN TERMS OF

THE NOW REMAINING N-2 VARIABLES IS PERFORMED. ONE CONTINUES
 IN THIS FASHION, ELIMINATION, ONE VARIABLE PER EQUATION,
 UNTIL FOR THE NTH EQUATION, WE ARE LEFT WITH ONE EQUATION
 IN ONE UNKNOWN. A SINGLE NEWTON STEP IS NOW PERFORMED, FOLLOWED BY BACK-SUBSTITUTION IN THE TRIANGULARIZED LINEAR SYSTEM GENERATING FOR THE X(I)'S. A PIVOTING THAT VARIABLE IS ACHIEVED BY CHOOSING FOR ELIMINATION AT ANY STEP THAT VARIABLE HAVING A PARTIAL DERIVATIVE OF LARGEST ABSOLUTE VALUE. THE PIVOTING IS DONE WITHOUT PHYSICAL INTERCHANGE OF ROWS OR COLUMNS.

REFERENCES

1. BROWN, K. M. "ALGORITHM 316 SOLUTION OF SIMULTANEOUS NON-LINEAR EQUATIONS (C5)", COLLECTED ALGORITHMS FROM CACM, JULY 1967, 316-P 1-0 AND 1-6.
 THE REMARKS IMMEDIATELY ABOVE ARE REPRODUCED FROM THIS WRITE-UP.
2. BROWN, K. M., A QUADRATICALLY CONVERGENT METHOD FOR SOLVING SIMULTANEOUS NON-LINEAR EQUATIONS. DOCTORAL THESIS, DEPT. COMPUTER SCIENCES, PURDUE U., LAFAYETTE, IND., AUG., 1966.
3. BROWN, K. M., AND CONF. S. D. THE SOLUTION OF SIMULTANEOUS NONLINEAR EQUATIONS. PROC. ACM 22ND NAT. CONF., PP 111-114.

.....

```

DIMENSION X(1)
DIMENSION U(13,13)
INTEGER CONVRG,TALLY,PONTER
COMMON PONTER(30,30),COE(30,30),ISUB(30),U,MM,NN
EXTERNAL EVALUT
DIMENSION TEMP(30),PART(30)

```

```

CONVRG=1
INITIALIZE CONVERGENCE AND TOLERANCE DETERMINERS.
ISING=1
RELCON=10.**(-NUMSIG)
ITERATION LOOP.
DO 55 M=1,MAXIT
  PROGRAM WILL PRINT OUT SUCCESSIVE APPROXIMATIONS OF X IF
  IPRINT IS 0.
  IF (IPRINT.EQ.0) WRITE (6,1000) M,(X(I),I=1,N)
  1000 FORMAT(1X,12HON ITERATION,14,17HTHE X ESTIMATE IS,5F15.7/7F15.7/

```

CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC

C C C C C


```

1      7F15.07(///)
C      INITIALIZE FIRST ROW OF ARRAY WHICH ENABLES INDEXING TO PIVOT
C      THROUGH THE VARIABLES WITHOUT INTERCHANGING ROWS OR COLUMNS.
5      DO 5 J=1,N
C      ONE K FOR EACH EQUATION.
C      DO 40 K=1,N
C      IF K IS GREATER THAN 1 BACK-SUBSTITUTION IS NECESSARY.
C      IF (K.GT.1) CALL BAKSUB(K,N,X)
C      CALL EVALUT(X,K,F)
C      FACTOR=0.001
7      TALLY=0
C      DO 10 I=K,N
C      ITEMP=PONTER(K,I)
C      HOLD=X(ITEMP)
C      GET INCREMENT TO OBTAIN ITEMP,TH PARTIAL.
C      H=FACTOR*HOLD
C      IF (H.EQ.0.) H=0.001
C      X(ITEMP)=HOLD+H
C      IF K IS GREATER THAN 1, A NEW X(ITEMP) VALUE WILL AFFECT X'S
C      WHICH ARE EXPRESSED IN TERMS OF IT AS A RESULT OF PREVIOUS
C      EQUATION STEP.
C      IF (K.GT.1) CALL BAKSUB(K,N,X)
C      CALL EVALUT(X,K,FPLUS)
C      GET ITEMP,TH PARTIAL.
C      PART(ITEMP)=(FPLUS-F)/H
C      X(ITEMP)=HOLD
C      IF PARTIAL IS TOO SMALL, INCREASE TALLY.
C      IF (ABS(PART(ITEMP)).NE.0.) GO TO 80
C      TALLY=TALLY+1
C      IF (ABS(F/PART(ITEMP)).GT.1.0E20) GO TO 85
85     GO TO 10
C      CONTINUE
C      IF (TALLY,LE,(N-K)) GO TO 15
C      FACTOR=FACTOR*.10
C      IF SURFACE IS TOO FLAT, THE SINGULARITY INDICATOR IS SET TO
C      ZERO AND RETURN IS EXECUTED.
C      IF (FACTOR.GT.0.5) GO TO 65
C      GO TO 7
15     IF (K.LT.N) GO TO 20
C      IF LAST PARTIAL IS ZERO, A SINGULARITY IS INDICATED AND A
C      RETURN EXECUTED.
C      IF (ABS(PART(ITEMP)).EQ.0.) GO TO 65
C      COE(K,N+1)=0
C      KMAX=ITEMP
C      GO TO 40
C      KMAX=PONTER(K,K)
20     DERMAT=ABS(PART(KMAX))

```

NLNS1610
 NLNS1620
 NLNS1630
 NLNS1640
 NLNS1650
 NLNS1660
 NLNS1670
 NLNS1680
 NLNS1690
 NLNS1700
 NLNS1710
 NLNS1720
 NLNS1730
 NLNS1740
 NLNS1750
 NLNS1760
 NLNS1770
 NLNS1780
 NLNS1790
 NLNS1800
 NLNS1810
 NLNS1820
 NLNS1830
 NLNS1840
 NLNS1850
 NLNS1860
 NLNS1870
 NLNS1880

NLNS1910
 NLNS1920
 NLNS1930
 NLNS1940
 NLNS1950
 NLNS1960
 NLNS1970
 NLNS1980
 NLNS1990
 NLNS2000
 NLNS2010
 NLNS2020
 NLNS2030
 NLNS2040
 NLNS2050
 NLNS2060


```

C      KPLUS=K+1
C      GET INDEX FOR LARGEST PARTIAL IN K'TH EQUATION.
C      DO 30 I=KPLUS,N
C      JSUB=PONTER(K,I)
C      TEST=ABS(PART(JSUB))
C      IF (TEST.LT.DERMAX) GO TO 25
C      DERMAX=TEST
C      DEFINE PIVOT TO SWIVEL ABOUT THE VARIABLE WITH MAXIMUM PARTIAL
C      WHEN WE GET TO THE NEXT EQUATION.
C      PONTER(KPLUS,I)=KMAX
C      IF THIS PARTIAL IS BIGGER, WE HAVE A NEW MAXIMUM.
C      KMAX=JSUB
C      GO TO 30
C      25 PONTER(KPLUS,I)=JSUB
C      30 CONTINUE
C      IF THAT PARTIAL IS 0, INDICATE A SINGULARITY AND RETURN.
C      IF (ABS(PART(KMAX)).EQ.0) GO TO 65
C      JSUB=PONTER(KPLUS,I)
C      SAVE THESE CONSTANTS FOR FUTURE USE.
C      COE(K,JSUB)=-PART(JSUB)/PART(KMAX)
C      GET PART OF EXPRESSION FOR THE NEW X(KMAX) VALUE
C      35 COE(K,N+1)=COE(K,N+1)+PART(JSUB)*X(JSUB)
C      40 COE(K,N+1)=(COE(K,N+1)-F)/PART(KMAX)+X(KMAX)
C      IF N IS 1, WE HAVE OUR SOLUTION IN THE NEXT STEP WITHOUT ANY
C      BACK-SUBSTITUTING.
C      X(KMAX)=COE(N,N+1)
C      FOR N GREATER THAN 1, WE PERFORM A FINAL BACK-SUBSTITUTION TO
C      GET OUR NEW X-VECTOR.
C      IF (N.GT.1) CALL BAKSUB(N,N,X)
C      43 DO 43 I=1,N
C      IF (M.EQ.1) GO TO 50
C      IF (ABS((TEMP(I)-X(I))/X(I)).GT.RELCON) GO TO 45
C      CONTINUE
C      CONVRG=CONVRG+1
C      IF IT CONVERGES, RETURN WITH LAST VECTOR.
C      IF (CONVRG.GE.3) GO TO 60
C      GO TO 50
C      45 CONVRG=1
C      SAVE CURRENT X-VECTOR FOR TESTING WITH NEXT X-VECTOR.
C      50 DO 55 I=1,N
C      55 TEMP(I)=X(I)
C      IF M IS THE ITERATION LIMIT, RETURN,
C      GO TO 70
C      60 MAXIT=M

```

NLNS2070
 NLNS2080
 NLNS2090
 NLNS2100
 NLNS2110
 NLNS2120
 NLNS2130
 NLNS2140
 NLNS2150
 NLNS2160
 NLNS2170
 NLNS2180
 NLNS2190
 NLNS2200
 NLNS2210
 NLNS2220
 NLNS2230
 NLNS2240
 NLNS2250
 NLNS2260
 NLNS2270
 NLNS2280
 NLNS2290
 NLNS2300
 NLNS2310
 NLNS2320
 NLNS2330
 NLNS2340
 NLNS2350
 NLNS2360
 NLNS2370
 NLNS2380
 NLNS2390
 NLNS2400
 NLNS2410
 NLNS2420
 NLNS2430
 NLNS2440
 NLNS2450
 NLNS2460
 NLNS2470
 NLNS2480
 NLNS2490
 NLNS2500
 NLNS2510
 NLNS2520
 NLNS2530
 NLNS2540

NLNS25550
NLNS25560
NLNS25570
NLNS25580
RKSR25590

GO TO 70
65 ISING=0
70 RETURN
END

C

BKSR26000
BKSR26010
BKSR26020

SUBROUTINE BAKSUB(K,N,X)
THIS SUBROUTINE BACK-SUBSTITUTES OR UPDATES VARIABLES WHICH ARE
FUNCTIONS OF CURRENT X ENTRY VALUES.
DIMENSION U(13,13)
INTEGER PONTER
COMMON PONTER(30,30),COE(30,30),ISUR(30),U,MM,NN

C

BKSR26300
BKSR26400
BKSR26500
BKSR26570
BKSR26580
BKSR26600
BKSR26700
BKSR26710
BKSR26720
BKSR26730
BKSR26740
BKSR26750
BKSR26760
BKSR26770

DO 10 KMM=2,K
KM=K+2-KMM
KMAX=ISUB(KM-1)
X(KMAX)=0.
DO 5 J=KM,N
JSUR=PONTER(KM,J)
SEE (2) FOR THE EXPRESSION FOR X(KMAX).
5 X(KMAX)=X(KMAX)+COE(KM-1,JSUR)*X(JSUR)
10 X(KMAX)=X(KMAX)+COE(KM-1,N+1)
RETURN
END

C

//GO%ETO6EQ01 DD SYSOUT=A,DCB=BLKSIZE=133
//GO%SYSDUMP DD SYSOUT=A
//GO%SYSD D D *

254 13
40205455 16131509130316071813061427172309
4020 54492624740
4020 21131616150609021816061115212413
4021 555454274010
40265371 24141517121414082109041811181413
4026 524427225144
40334568 16061413141409072406122321131811
4033 4742230294944
40504147 17102205160906142510111621171114
4050 475137314024
4050 23111709110816032115121623132012
4050 455424334129
40625774 25121806120610231104081823191511
4062 494038263048
40694247 19101514180000072308082219281614
4069 475733264829
40705675 19152221201007062112011619041613
4070 505034234637

(Sample data for
9 points)

VI. REFERENCES

1. Bonner, R. E., "On Some Clustering Techniques," IBM Journal of Research Development, v. 8, p. 22, 1964.
2. Dixon, J. W., Biomedical Computer Programs, Berkeley, pp. 196-206, 1968.
3. Dupraw, E. J., "Non-linnean Taxonomy," Nature, v. 202, pp. 849-852, May 30, 1964.
4. Edwards, A. W. F. and Cavalli-Sforza, "A Method for Cluster Analysis," Biometrics, v. 21, p. 62, 1964.
5. Gilmore, J. S. L., "The Development of Taxonomic Theory Since 1851," Nature, v. 168, pp. 400-402, 1951.
6. Gilmore, J. S. L., "A Taxonomic Problem," Nature, v. 140, pp. 1040-1042, 1937.
7. Gilmore, J. S. L., The New Systematics, Claredon, Oxford, 1940.
8. Gower, J. C., "A Comparison of Some Methods of Cluster Analysis," Biometrics, v. 23, pp. 623-637, Dec., 1967.
9. Hadley, G., Nonlinear and Dynamic Programming, Addison-Wesley, Palo Alto, 1964.
10. Hotelling, H., "The Relations Between Two Sets of Variates," Biometrika, v. 28, p. 328, 1936.
11. Rao, C. R., Advanced Statistical Methods in Biometric Research, Wiley, New York, 1952.
12. Rogers, D. J. and Tanimoto, T. T., "A Computer Program for Classifying Plants," Science, v. 132, pp. 1115-1118, 1960.
13. Simpson, G. G., Principles of Animal Taxonomy, Columbia, New York, 1961.
14. Sneath, P. H. A., and Sokal, R. R., "Numerical Taxonomy," Nature, v. 193, pp. 855-860, 1962.
15. Sokal, R. R., "Numerical Taxonomy," Science American, v. 215, pp. 106-116, Dec., 1966.

16. Sokal, R. R. and Sneath, P. H. A., Principles of Numerical Taxonomy, W. H. Freeman Co., San Francisco, 1963.
17. Spranger, E., Types of Men, Stechert-Hafner, Inc., New York

INITIAL DISTRIBUTION LIST

	<u>No. Copies</u>
1. Defense Documentation Center Cameron Station Alexandria, Virginia 22314	2
2. Library, Code 0212 Naval Postgraduate School Monterey, California 93940	2
3. Asst. Prof. G. A. Tuck, Code 55 Tk (Adviser) Department of Operations Analysis Naval Postgraduate School Monterey, California 93940	1
4. LT (junior grade) William M. Cima 344 Ridge Avenue Pittsburgh, Pennsylvania 15221	1
5. Prof. J. R. Borsting, Code 55 Bg Department of Operations Analysis Naval Postgraduate School Monterey, California 93940	1

DOCUMENT CONTROL DATA - R & D

(Security classification of title, body of abstract and indexing annotation must be entered when the overall report is classified)

1. ORIGINATING ACTIVITY (Corporate author) Naval Postgraduate School Monterey, California 93940		2a. REPORT SECURITY CLASSIFICATION Unclassified	
		2b. GROUP	
3. REPORT TITLE A Method of Cluster Analysis			
4. DESCRIPTIVE NOTES (Type of report and inclusive dates) Master's Thesis; June 1970			
5. AUTHOR(S) (First name, middle initial, last name) William M. Cima			
6. REPORT DATE June 1970		7a. TOTAL NO. OF PAGES 76	7b. NO. OF REFS 17
8a. CONTRACT OR GRANT NO.		9a. ORIGINATOR'S REPORT NUMBER(S)	
b. PROJECT NO.			
c.		9b. OTHER REPORT NO(S) (Any other numbers that may be assigned this report)	
d.			
10. DISTRIBUTION STATEMENT This document has been approved for public release and sale; its distribution is unlimited.			
11. SUPPLEMENTARY NOTES		12. SPONSORING MILITARY ACTIVITY Naval Postgraduate School Monterey, California 93940	
13. ABSTRACT A method of cluster analysis is presented in which points in n-dimensional space are analyzed through a subdivisive procedure. The points are orthogonally projected onto that line which maximizes their variance and the resulting point distribution is then analyzed with the use of a histogram. Wherever possible, divisions between conglomerates of points are made and each separate clump is subsequently analyzed. Ultimately adjacent groups are combined and analyzed through an analogous technique in an effort to re-unite any points which may have inadvertently deviated from the group with which they truly associate. The method is later refined to allow the detection of groups in several point dispersions which would have appeared as a single conglomeration under the original method. An example is given to illustrate the applicability of the procedure.			

14

KEY WORDS

LINK A

LINK B

LINK C

ROLE

WT

ROLE

WT

ROLE

WT

Maximum variance line

Relative maximum variance line

Conglomerates

Orthogonal projection

Cluster analysis

Statistics

Thesis
C4789
c.2

Cima

A method of cluster
analysis.

121707

20 APR 73
S SEP 78
11 SEP 79
25 APR 87

21149
~~24807~~
26031
31585

Thesis
C4789
c.2

Cima

A method of cluster
analysis.

121707

thesC4789
A method of cluster analysis.



3 2768 001 02671 9
DUDLEY KNOX LIBRARY